



# Identification of polynomial chaos representations in high dimension from a set of realizations

Guillaume Perrin, Christian Soize, Denis Duhamel, Christine Fünfschilling

## ► To cite this version:

Guillaume Perrin, Christian Soize, Denis Duhamel, Christine Fünfschilling. Identification of polynomial chaos representations in high dimension from a set of realizations. SIAM Journal on Scientific Computing, 2012, 34 (6), pp.A2917-A2945. 10.1137/11084950X . hal-00770006

**HAL Id: hal-00770006**

**<https://hal.science/hal-00770006>**

Submitted on 4 Jan 2013

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# IDENTIFICATION OF POLYNOMIAL CHAOS REPRESENTATIONS IN HIGH DIMENSION FROM A SET OF REALIZATIONS

G. PERRIN<sup>\*†‡</sup>, C. SOIZE<sup>\*</sup>, D. DUHAMEL<sup>†</sup>, AND C. FUNFSCHILLING<sup>‡</sup>

## Abstract.

This paper deals with the identification in high dimension of polynomial chaos expansion of random vectors from a set of realizations. Due to numerical and memory constraints, the usual polynomial chaos identification methods are based on a series of truncations that induces a numerical bias. This bias becomes very detrimental to the convergence analysis of polynomial chaos identification in high dimension. This paper therefore proposes a new formulation of the usual polynomial chaos identification algorithms to avoid this numerical bias. After a review of the polynomial chaos identification method, the influence of the numerical bias on the identification accuracy is quantified. The new formulation is then described in details, and illustrated on two examples.

**Key words.** polynomial chaos expansion, high dimension, computation.

**AMS subject classifications.** 60H35, 60H15, 60H25, 60H40, 65C50

**1. Introduction.** In spite of always more accurate numerical solvers, deterministic models are not able to represent most of the experimental data, which are variable and often uncertain by nature. Hence, the application fields of non deterministic modeling, which can take into account the model parameters variability as well as the model error uncertainties, has kept increasing. Uncertainties are therefore introduced in computational mechanical models with more and more degrees of freedom. In this context, the characterization of the probability distribution  $P_{\boldsymbol{\eta}}(d\mathbf{x})$  of  $N_{\boldsymbol{\eta}}$ -dimension random vector  $\boldsymbol{\eta}$  from sets of experimental measurements is bound to play a key role, in particular, in high dimension, that is to say for a large value of  $N_{\boldsymbol{\eta}}$ . In this work, it is assumed that  $P_{\boldsymbol{\eta}}(d\mathbf{x}) = p_{\boldsymbol{\eta}}(\mathbf{x})d\mathbf{x}$  in which the probability density function (PDF)  $p_{\boldsymbol{\eta}}$  is a function in the set  $\mathcal{F}(\mathcal{D}, \mathbb{R}^+)$  of all the positive-valued functions defined on any part  $\mathcal{D}$  of  $\mathbb{R}^{N_{\boldsymbol{\eta}}}$  and for which integral over  $\mathcal{D}$  is 1.

Two kinds of methods can be used to build such a PDF: the direct and the indirect methods. Among the direct methods, the Prior Algebraic Stochastic Modeling (PASM) methods postulate an algebraic representation  $\boldsymbol{\eta} \approx t^{\text{alg}}(\boldsymbol{\Xi}, \mathbf{w})$ , with  $t^{\text{alg}}$  a prior transformation,  $\boldsymbol{\Xi}$  a given random vector and  $\mathbf{w}$  a vector of parameters to identify. In the same category, the methods based on the Information Theory and the Maximum Entropy Principle (MEP) have been developped (see [13] and [27]) to compute  $p_{\boldsymbol{\eta}}$  from the only available information of random vector  $\boldsymbol{\eta}$ . This information can be seen as the admissible set  $\mathcal{C}^{\text{ad}}$  for  $p_{\boldsymbol{\eta}}$ :

$$\mathcal{C}^{\text{ad}} = \left\{ p_{\boldsymbol{\eta}} \in \mathcal{F}(\mathcal{D}, \mathbb{R}^+) \mid \int_{\mathcal{D}} p_{\boldsymbol{\eta}}(\mathbf{x})d\mathbf{x} = 1, \right. \\ \left. \forall 1 \leq m \leq M, \int_{\mathcal{D}} \mathbf{g}_m(\mathbf{x})p_{\boldsymbol{\eta}}(\mathbf{x})d\mathbf{x} = \mathbf{f}_m \right\}, \quad (1.1)$$

<sup>\*</sup>Université Paris-Est, Modélisation et Simulation Multi-Échelle (MSME UMR 8208 CNRS), 5 Bd. Descartes, 77454 Marne-la-Vallée, France (christian.soize@univ-paris-est.fr).

<sup>†</sup>Université Paris-Est, Navier (ENPC-IFSTTAR-CNRS UMR 8205), Ecole Nationale des Ponts et Chaussées, 6 et 8 Avenue Blaise Pascal, Cité Descartes, Champs sur Marne, 77455 Marne-la-Vallée, Cedex 2, France (denis.duhamel@enpc.fr)

<sup>‡</sup>SNCF, Innovation and Research Department, Immeuble Lumière, 40 avenue des Terroirs de France, 75611, Paris, Cedex 12, France (guillaume.perrin@sncf.fr, christine.funfschilling@sncf.fr).

where  $\{\mathbf{f}_m, 1 \leq m \leq M\}$  gathers  $M$  given vectors which are respectively associated with a given vector-valued functions  $\{\mathbf{g}_m, 1 \leq m \leq M\}$ . Hence, the MPE allows building  $p_{\boldsymbol{\eta}}$  as the solution of the optimization problem:

$$p_{\boldsymbol{\eta}} = \arg \max_{p_{\boldsymbol{\eta}} \in \mathcal{C}^{\text{ad}}} \left\{ - \int_{\mathcal{D}} p_{\boldsymbol{\eta}}(\mathbf{x}) \log(p_{\boldsymbol{\eta}}(\mathbf{x})) d\mathbf{x} \right\}. \quad (1.2)$$

On the other hand, the indirect methods allow the construction of the PDF  $p_{\boldsymbol{\eta}}$  of the considered random vector  $\boldsymbol{\eta}$  from a transformation  $\mathbf{t}$  of a known random vector  $\boldsymbol{\xi} = (\xi_1, \dots, \xi_{N_g})$  of given dimension  $N_g \leq N_{\boldsymbol{\eta}}$ :

$$\boldsymbol{\eta} = \mathbf{t}(\boldsymbol{\xi}), \quad (1.3)$$

defining a transformation  $\mathbb{T}$  between  $p_{\boldsymbol{\eta}}$  and the PDF  $p_{\boldsymbol{\xi}}$  of  $\boldsymbol{\xi}$ :

$$p_{\boldsymbol{\eta}} = \mathbb{T}(p_{\boldsymbol{\xi}}). \quad (1.4)$$

The construction of the transformation  $\mathbf{t}$  is thus the key point of these indirect methods. In this context, the isoprobabilistic transformations such as the Nataf transformation (see [20]) or the Rosenblatt transformation (see [23]) have allowed the development of interesting results in the second part of the twentieth century but are still limited to very small dimension cases and not to the high dimension case considered in this work. Nowadays, the most popular indirect methods are the polynomial chaos expansion (PCE) methods, which have been first introduced by Wiener [33] for stochastic processes, and pioneered by Ghanem and Spanos [10] [11] for the use of it in computational sciences. In the last decade, this very promising method has thus been applied in many works (see, for instance [1], [2], [3], [4], [5], [7], [8], [9], [12], [14], [15], [16], [19], [18], [17], [21], [22], [24], [26], [28], [31], [32], [25], [34]). The PCE is based on a direct projection of the random vector  $\boldsymbol{\eta}$  on a chosen hilbertian basis  $\mathcal{B}_{\text{orth}} = \{\psi_{\boldsymbol{\alpha}}(\boldsymbol{\xi}), \boldsymbol{\alpha} \in \mathbb{N}^{N_g}\}$  of all the second-order random vectors with values in  $\mathbb{R}^{N_{\boldsymbol{\eta}}}$ :

$$\boldsymbol{\eta} = \sum_{\boldsymbol{\alpha} \in \mathbb{N}^{N_g}} \mathbf{y}^{(\boldsymbol{\alpha})} \psi_{\boldsymbol{\alpha}}(\boldsymbol{\xi}), \quad (1.5)$$

$$\boldsymbol{\xi} \mapsto \psi_{\boldsymbol{\alpha}}(\boldsymbol{\xi}) = X_{\alpha_1}(\xi_1) \otimes \dots \otimes X_{\alpha_{N_g}}(\xi_{N_g}), \quad (1.6)$$

where  $x \mapsto X_{\alpha_{\ell}}(x)$  is the normalized polynomial basis of degree  $\alpha_{\ell}$  associated with the PDF  $p_{\xi_{\ell}}$  of the random variable  $\xi_{\ell}$ , and  $\boldsymbol{\alpha}$  is the multi-index of the multidimensional polynomial basis element  $\psi_{\boldsymbol{\alpha}}(\boldsymbol{\xi})$ . Building the transformation  $\mathbf{t}$  requires therefore the construction of the projection vectors  $\{\mathbf{y}^{(\boldsymbol{\alpha})}, \boldsymbol{\alpha} \in \mathbb{N}^{N_g}\}$ .

The present work is devoted to the identification in high dimension of the PCE coefficients  $\{\mathbf{y}^{(\boldsymbol{\alpha})}, \boldsymbol{\alpha} \in \mathbb{N}^{N_g}\}$ , when the only available information on the random vector  $\boldsymbol{\eta}$  is a set of  $\nu^{\text{exp}}$  independent realizations  $\{\boldsymbol{\eta}^{(1)}, \dots, \boldsymbol{\eta}^{(\nu^{\text{exp}})}\}$ .

In practice, the PCE of  $\boldsymbol{\eta}$  has first to be truncated:

$$\boldsymbol{\eta} \approx \boldsymbol{\eta}^{\text{chaos}}(N) = \sum_{\boldsymbol{\alpha} \in \mathcal{A}_p} \mathbf{y}^{(\boldsymbol{\alpha})} \psi_{\boldsymbol{\alpha}}(\boldsymbol{\xi}), \quad (1.7)$$

$$\mathcal{A}_p = \left\{ \boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_{N_g}) \mid |\boldsymbol{\alpha}| = \sum_{\ell=1}^{N_g} \alpha_{\ell} \leq p \right\} = \left\{ \boldsymbol{\alpha}^{(1)}, \dots, \boldsymbol{\alpha}^{(N)} \right\}, \quad (1.8)$$

where  $\boldsymbol{\eta}^{\text{chaos}}(N)$  is the projection of  $\boldsymbol{\eta}$  on the  $N$ -dimension subspace spanned by  $\{ \psi_{\boldsymbol{\alpha}}(\boldsymbol{\xi}), \boldsymbol{\alpha} \in \mathcal{A}_p \} \subset \mathcal{B}_{\text{orth}}$ . It can be noticed that  $N$  increases very quickly with respect to the dimension  $N_g$  of  $\boldsymbol{\xi}$  and the maximum degree  $p$  of the truncated basis  $\{ \psi_{\boldsymbol{\alpha}}(\boldsymbol{\xi}), \boldsymbol{\alpha} \in \mathcal{A}_p \}$ , as:

$$N = (N_g + p)! / (N_g! p!). \quad (1.9)$$

Methods to perform the convergence analysis in high dimension with respect to a given error threshold on the PCE residue  $\boldsymbol{\eta} - \boldsymbol{\eta}^{\text{chaos}}(N)$  are therefore of great concern to justify the truncation parameters  $N_g$  and  $p$ .

In this prospect, the article [29] provides advanced algorithms to compute the PCE coefficients from the  $\nu^{\text{exp}}$  independent realizations  $\{ \boldsymbol{\eta}^{(1)}, \dots, \boldsymbol{\eta}^{(\nu^{\text{exp}})} \}$  by focusing on the maximization of the likelihood. In particular, one of the key point of these algorithms is the calculation of  $(N \times \nu^{\text{chaos}})$  real matrix  $[\Psi]$  of independent realizations of the truncated PCE basis  $\{ \psi_{\boldsymbol{\alpha}}(\boldsymbol{\xi}), \boldsymbol{\alpha} \in \mathcal{A}_p \}$ :

$$[\Psi] = [\Psi(\boldsymbol{\xi}(\theta_1), p) \quad \dots \quad \Psi(\boldsymbol{\xi}(\theta_{\nu^{\text{chaos}}}), p)], \quad (1.10)$$

$$\Psi(\boldsymbol{\xi}, p) = (\psi_{\boldsymbol{\alpha}^{(1)}}(\xi_1, \dots, \xi_{N_g}), \dots, \psi_{\boldsymbol{\alpha}^{(N)}}(\xi_1, \dots, \xi_{N_g})), \quad (1.11)$$

where the set  $\{ \boldsymbol{\xi}(\theta_1), \dots, \boldsymbol{\xi}(\theta_{\nu^{\text{chaos}}}) \}$  gathers  $\nu^{\text{chaos}}$  independent realizations of the random vector  $\boldsymbol{\xi}$ .

Recurrence formula or algebraic explicit representations are generally used to compute such matrix  $[\Psi]$ , which are supposed to verify the asymptotical property:

$$\lim_{\nu^{\text{chaos}} \rightarrow +\infty} \frac{1}{\nu^{\text{chaos}}} [\Psi][\Psi]^T = [I_N], \quad (1.12)$$

as a direct consequence of the orthonormality of the PCE basis  $\{ \psi_{\boldsymbol{\alpha}}, \boldsymbol{\alpha} \in \mathcal{A}_p \}$ , where  $[I_N]$  is the  $N$ -dimension identity matrix.

However, for numerically admissible values of  $\nu^{\text{chaos}}$  (between 1000 and 10000), it has been shown in [30] that the difference  $\frac{1}{\nu^{\text{chaos}}} [\Psi][\Psi]^T - [I_N]$  can be very significant when high values of the maximum degree  $p$  can be encountered with simultaneously significant values of  $N_g$ . This difference induces a detrimental bias in the PCE identification, which makes the convergence of classical PCE in high dimension very difficult.

In [30], it is therefore proposed a method using singular matrix decomposition to numerically adapt classical generations of  $[\Psi]$ , and make this difference be zero for any values of  $p$  and  $N_g$ . Nevertheless, this conditionning on  $[\Psi]$  modifies the initial structure of  $[\Psi]$ , and makes the identified PCE coefficients  $\{\mathbf{y}^{(\alpha)}, \alpha \in \mathcal{A}_p\}$  impossible to be reused on an other matrix  $[\Psi^*]$  of  $\nu^{\text{chaos},*}$  new realizations of  $\Psi(\xi, p)$ .

As an extension of the works described in [29] and [30], this article proposes an original decomposition of the PCE coefficients  $\{\mathbf{y}^{(\alpha)}, \alpha \in \mathcal{A}_p\}$ , that reduces the numerical bias introduced during the identification by the finite dimension of  $[\Psi]$  and for large values of degree  $p$ . This new formulation is particularly adapted to the high dimension, and allows the identified coefficients to be reused for other matrix of realizations  $[\Psi^*]$ .

In Section 2, the PCE identification from a set of experimental data with an arbitrary measure is described. In particular, the role played by the matrix of independent realizations  $[\Psi]$  is emphasized. Section 3 focuses on the convergence properties of this matrix  $[\Psi]$  with respect to three statistical measures, and describes an innovative method to generate this matrix without using computational recurrence formula nor algebraic explicit representation. In Section 4, the new formulation of the PCE identification problem is given. Finally, are presented in Section 5 two applications of the former method with a Gaussian measure.

**2. PCE identification of random vectors from a set of independent realizations.** In this section, a description of the PCE identification with respect to an arbitrary measure is given. The objective is to summarize the different key steps of the PCE identification method and the way they are practically implemented.

After having defined the theoretical frame of the PCE identification, the cost-function that leads to the computation of the PCE coefficients  $\{\mathbf{y}^{(\alpha)}, \alpha \in \mathcal{A}_p\}$  is presented, for given truncation parameters  $N_g$  and  $p$ . At last, to justify the choice of these truncation parameters, a method to perform the convergence analysis is introduced.

**2.1. Theoretical frame.** Let  $(\Theta, \mathcal{T}, \mathcal{P})$  be a probability space. Let  $L_{\mathcal{P}}^2(\Theta, \mathbb{R}^{N_\eta})$  be the space of all the second-order  $N_\eta$ -dimension random vectors defined on  $(\Theta, \mathcal{T}, \mathcal{P})$  with values in  $\mathbb{R}^{N_\eta}$ , equipped with the inner product  $\langle \cdot, \cdot \rangle$ :

$$\langle \mathbf{U}, \mathbf{V} \rangle = \int_{\Theta} \mathbf{U}^T(\theta) \mathbf{V}(\theta) dP(\theta) = E(\mathbf{U}^T \mathbf{V}), \quad \forall \mathbf{U}, \mathbf{V} \in L_{\mathcal{P}}^2(\Theta, \mathbb{R}^{N_\eta}), \quad (2.1)$$

where  $E(\cdot)$  is the mathematical expectation.

Let  $\boldsymbol{\eta} = (\eta_1, \dots, \eta_{N_\eta})$  be an element of  $L_{\mathcal{P}}^2(\Theta, \mathbb{R}^{N_\eta})$ . It is assumed that  $\nu^{\text{exp}}$  independent realizations  $\{\boldsymbol{\eta}^{(1)}, \dots, \boldsymbol{\eta}^{(\nu^{\text{exp}})}\}$  of  $\boldsymbol{\eta}$  are known and gathered in the  $(N_\eta \times \nu^{\text{exp}})$  real matrix  $[\boldsymbol{\eta}^{\text{exp}}]$ :

$$[\boldsymbol{\eta}^{\text{exp}}] = \begin{bmatrix} \boldsymbol{\eta}^{(1)} & \dots & \boldsymbol{\eta}^{(\nu^{\text{exp}})} \end{bmatrix}. \quad (2.2)$$

Equation (1.7) can be rewritten as:

$$\boldsymbol{\eta}^{\text{chaos}}(N) = [\mathbf{y}] \Psi(\xi, p), \quad (2.3)$$

$$[y] = \begin{bmatrix} \mathbf{y}^{(\alpha^{(1)})} & \dots & \mathbf{y}^{(\alpha^{(N)})} \end{bmatrix}. \quad (2.4)$$

The orthonormality property of the projection basis  $\{\psi_{\alpha}(\xi), \alpha \in \mathcal{A}_p\}$  yields the condition:

$$E(\Psi(\xi, p)\Psi(\xi, p)^T) = [I_N]. \quad (2.5)$$

Since  $\psi_{\alpha^{(1)}}(\xi) = 1$ , it can be seen that:

$$E(\boldsymbol{\eta}^{\text{chaos}}(N)) = \mathbf{y}^{(\alpha^{(1)})}. \quad (2.6)$$

Let  $[R_{\eta}]$  and  $[R_{\eta}^{\text{chaos}}(N)]$  be the autocorrelation matrix of the random vectors  $\boldsymbol{\eta}$  and  $\boldsymbol{\eta}^{\text{chaos}}(N)$ :

$$[R_{\eta}] = E(\boldsymbol{\eta}\boldsymbol{\eta}^T), \quad (2.7)$$

$$[R_{\eta}^{\text{chaos}}(N)] = E\left(\boldsymbol{\eta}^{\text{chaos}}(N)(\boldsymbol{\eta}^{\text{chaos}}(N))^T\right) = [y]E(\Psi(\xi, p)\Psi(\xi, p)^T)[y]^T = [y][y]^T. \quad (2.8)$$

**2.2. Identification of the polynomial chaos expansion coefficients.** In this section, particular values of the truncation parameters  $N_g$  and  $p$  are considered. Let  $\mathcal{M}_{N_{\eta}N}$  be the space of all the  $(N_{\eta} \times N)$  real matrices. For a given value of  $[y^*]$  in  $\mathcal{M}_{N_{\eta}N}$ , the random vector  $\mathbf{U}([y^*]) = [y^*]\Psi(\xi, p)$  is a  $N_{\eta}$ -dimension random vector, for which the autocorrelation is equal to  $[y^*][y^*]^T$ . Let  $p_{\mathbf{U}([y^*])}$  be its multidimensional PDF.

When the only available information on  $\boldsymbol{\eta}$  is a set of  $\nu^{\text{exp}}$  independent realizations, the optimal coefficients matrix  $[y]$  of its truncated PCE,  $\boldsymbol{\eta}^{\text{chaos}}(N) = [y]\Psi(\xi, p)$ , can be seen as the argument which maximizes the log-likelihood  $\mathcal{L}_{\mathbf{U}([y^*])}([\eta^{\text{exp}}])$  of  $\mathbf{U}([y^*])$ :

$$[y] = \arg \max_{[y^*] \in \mathcal{M}_{N_{\eta}N}} \mathcal{L}_{\mathbf{U}([y^*])}([\eta^{\text{exp}}]), \quad (2.9)$$

$$\mathcal{L}_{\mathbf{U}([y^*])}([\eta^{\text{exp}}]) = \sum_{i=1}^{\nu^{\text{exp}}} \ln p_{\mathbf{U}([y^*])}(\boldsymbol{\eta}^{(i)}). \quad (2.10)$$

### 2.3. Practical solving of the log-likelihood maximization.

#### 2.3.1. The need for statistical algorithms to maximize the log-likelihood.

The log-likelihood  $\mathcal{L}_{\mathbf{U}([y^*])}([\eta^{\text{exp}}])$  being non-convex, deterministic algorithms such as gradient algorithms cannot be applied to solve Eq. (2.9), and random search algorithms have to be used. Hence, the precision of the PCE has to be correlated to a numerical cost  $M$ , which corresponds to a number of independent trials of  $[y^*]$  in  $\mathcal{M}_{N_{\eta}N}$ . Let  $\mathcal{Y} = \{[y^*]^{(r)}, 1 \leq r \leq M\}$  be a set of  $M$  elements, which have been

chosen randomly in  $\mathcal{M}_{N_\eta N}$ . For a given numerical cost  $M$ , the most accurate PCE coefficients matrix  $[y]$  is approximated by:

$$[y] \approx [y_{\mathcal{Y}}] = \arg \max_{[y^*] \in \mathcal{Y}} \mathcal{L}_{\mathcal{U}([y^*])}([\eta^{\text{exp}}]). \quad (2.11)$$

**2.3.2. Restriction of the maximization domain.** From the  $\nu^{\text{exp}}$  independent realizations  $\{\boldsymbol{\eta}^{(1)}, \dots, \boldsymbol{\eta}^{(\nu^{\text{exp}})}\}$ , the mean value  $E(\boldsymbol{\eta})$  and the autocorrelation matrix  $[R_\eta]$  of  $\boldsymbol{\eta}$  can be estimated by:

$$E(\boldsymbol{\eta}) \approx \hat{\boldsymbol{\eta}}(\nu^{\text{exp}}) = \frac{1}{\nu^{\text{exp}}} \sum_{i=1}^{\nu^{\text{exp}}} \boldsymbol{\eta}^{(i)}, \quad (2.12)$$

$$[R_\eta] \approx [\hat{R}_\eta(\nu^{\text{exp}})] = \frac{1}{\nu^{\text{exp}}} \sum_{i=1}^{\nu^{\text{exp}}} \boldsymbol{\eta}^{(i)} (\boldsymbol{\eta}^{(i)})^T = \frac{1}{\nu^{\text{exp}}} [\eta^{\text{exp}}][\eta^{\text{exp}}]^T. \quad (2.13)$$

A good way to improve the efficiency of the numerical identification of  $[y]$  is then to restrict the research set to  $\mathcal{O}_\eta \subset \mathcal{M}_{N_\eta N}$ , with:

$$\begin{aligned} \mathcal{O}_\eta = \left\{ [y] = [\mathbf{y}^{(\alpha^{(1)})}, \dots, \mathbf{y}^{(\alpha^{(N)})}] \in \mathcal{M}_{N_\eta N} \mid \right. \\ \left. \mathbf{y}^{(\alpha^{(1)})} = \hat{\boldsymbol{\eta}}(\nu^{\text{exp}}), [y][y]^T = [\hat{R}_\eta(\nu^{\text{exp}})] \right\}, \end{aligned} \quad (2.14)$$

which, taking into account Eqs. (2.6) and (2.8), guarantees by construction that:

$$\begin{cases} [R_\eta^{\text{chaos}}(N)] = [\hat{R}_\eta(\nu^{\text{exp}})], \\ E(\boldsymbol{\eta}^{\text{chaos}}(N)) = \hat{\boldsymbol{\eta}}(\nu^{\text{exp}}). \end{cases} \quad (2.15)$$

Hence, the PCE coefficients matrix  $[y]$  can be approximated as the argument in  $\mathcal{O}_\eta$  that maximizes the log-likelihood  $\mathcal{L}_{\mathcal{U}([y^*])}([\eta^{\text{exp}}])$ . By defining  $\mathcal{W}$  the set that gathers  $M$  randomly raised elements of  $\mathcal{O}_\eta$ ,  $[y]$  can then be assessed as the solution of the new optimization problem:

$$[y] \approx [y_{\mathcal{W}}] = \arg \max_{[y^*] \in \mathcal{W}} \mathcal{L}_{\mathcal{U}([y^*])}([\eta^{\text{exp}}]). \quad (2.16)$$

**2.3.3. Approximation of the log-likelihood function.** From a particular matrix of realizations  $[\Psi]$  (which is defined in Eq. (1.10)), if  $[y^*]$  is an element of  $\mathcal{O}_\eta$ ,  $\nu^{\text{chaos}}$  independent realizations  $\{\mathbf{U}([y^*], \theta_n) = [y^*]\boldsymbol{\Psi}(\boldsymbol{\xi}(\theta_n), p), 1 \leq n \leq \nu^{\text{chaos}}\}$  of the random vector  $\mathbf{U}([y^*])$  can be computed and gathered in the matrix  $[U]$ :

$$[U] = [\mathbf{U}([y^*], \theta_1) \quad \dots \quad \mathbf{U}([y^*], \theta_{\nu^{\text{chaos}}})] = [y^*][\Psi]. \quad (2.17)$$

Hence, using Gaussian Kernels, the PDF  $p_{\mathbf{U}([y^*])}$  of  $\mathbf{U}([y^*])$  can be directly estimated by its non parametric estimator  $\hat{p}_{\mathbf{U}}$ :

$$\forall \mathbf{x} \in \mathbb{R}^{N_\eta}, \quad p_{\mathbf{U}([y^*])}(\mathbf{x}) \approx \hat{p}_{\mathbf{U}}(\mathbf{x}) = \frac{1}{(2\pi)^{N_\eta/2} \nu^{\text{chaos}} \prod_{k=1}^{N_\eta} h_k} \sum_{n=1}^{\nu^{\text{chaos}}} \exp \left( -\frac{1}{2} \sum_{k=1}^{N_\eta} \left( \frac{x_k - U_k([y^*], \theta_n)}{h_k} \right)^2 \right), \quad (2.18)$$

where  $\mathbf{h} = (h_1, \dots, h_{N_\eta})$  is the multidimensionnal optimal Silverman bandwidth vector (see [6]) of the Kernel smoothing estimation of  $p_{\mathbf{U}([y^*])}$ :

$$\forall 1 \leq k \leq N_\eta, \quad h_k = \hat{\sigma}_{U_k} \left( \frac{4}{(2 + N_\eta) \nu^{\text{exp}}} \right)^{1/(N_\eta+4)}, \quad (2.19)$$

where  $\hat{\sigma}_{U_k}$  is the empirical estimation of the standard deviation of each component  $U_k$  of  $\mathbf{U}$ . It has to be noticed that  $\hat{p}_{\mathbf{U}}$  only depends on the bandwidth vector  $\mathbf{h}$ , and the two matrices  $[y^*]$  and  $[\Psi]$ . Hence, according to the Eqs. (2.10), (2.17) and (2.18), for a given value of  $\nu^{\text{chaos}}$ , the maximization of the log-likelihood function  $\mathcal{L}_{\mathbf{U}([y^*])}$  can be replaced by the maximization of the cost-function  $\mathcal{C}([\eta^{\text{exp}}], [y^*], [\Psi])$  such that:

$$[y] \approx [y_{\mathcal{O}_\eta}] = \arg \max_{[y^*] \in \mathcal{O}_\eta} \mathcal{C}([\eta^{\text{exp}}], [y^*], [\Psi]), \quad (2.20)$$

where:

$$\mathcal{C}([\eta^{\text{exp}}], [y^*], [\Psi]) = \mathcal{C}_C + \mathcal{C}_V([\eta^{\text{exp}}], [y^*], [\Psi]), \quad (2.21)$$

$$\mathcal{C}_C = -\nu^{\text{exp}} \ln \left( (2\pi)^{N_\eta/2} \nu^{\text{chaos}} \prod_{k=1}^{N_\eta} h_k \right), \quad (2.22)$$

$$\mathcal{C}_V([\eta^{\text{exp}}], [y^*], [\Psi]) = \sum_{i=1}^{\nu^{\text{exp}}} \ln \left( \sum_{n=1}^{\nu^{\text{chaos}}} \exp \left( -\frac{1}{2} \sum_{k=1}^{N_\eta} \left( \frac{\eta_k^{(i)} - U_k([y^*], \theta_n)}{h_k} \right)^2 \right) \right). \quad (2.23)$$

Hence, the optimization problem defined by Eq. (2.16) can finally be estimated by:

$$[y] \approx [y_{\mathcal{O}_\eta}^M] = \arg \max_{[y^*] \in \mathcal{W}} \mathcal{C}([\nu^{\text{exp}}], [y^*], [\Psi]). \quad (2.24)$$

The optimization problem defined by Eq. (2.24) is now supposed to be solved with the advanced algorithms described in [29] to optimize the trials of the elements of  $\mathcal{W}$  for a given computation cost  $M$ . The higher the value of  $M$  is, the better the PCE identification should be. Therefore, this value has to be chosen as high as possible while respecting the computational resource limitation.



**2.3.4. Accuracy of the PCE identification.** For a given computation cost  $M$ , let  $[y_{\mathcal{O}_\eta}^M]$  be an optimal solution of Eq. (2.24).  $[y_{\mathcal{O}_\eta}^M]$  is a numerical estimation of the PCE coefficients matrix  $[y]$ . For a new  $(N \times \nu^{\text{chaos},*})$  real matrix  $[\Psi^*]$  of independent realizations ( $\nu^{\text{chaos},*}$  can be higher than  $\nu^{\text{chaos}}$ ), the robustness of  $[y_{\mathcal{O}_\eta}^M]$  regarding the choice of  $[\Psi]$  can then be estimated by comparing  $\mathcal{C}([ \eta^{\text{exp}}, [y_{\mathcal{O}_\eta}^M], [\Psi] )$  and  $\mathcal{C}([ \eta^{\text{exp}}, [y_{\mathcal{O}_\eta}^M], [\Psi^*] )$ . In addition, if  $\nu^{\text{exp}}$  new independent realizations of  $\boldsymbol{\eta}$  were available and gathered in the matrix  $[\eta^{\text{exp,new}}]$ , the over-learning of the method could be measured by comparing  $\mathcal{C}([ \eta^{\text{exp}}, [y_{\mathcal{O}_\eta}^M], [\Psi] )$  and  $\mathcal{C}([ \eta^{\text{exp,new}}, [y_{\mathcal{O}_\eta}^M], [\Psi] )$ . At last, for the same computation cost  $M$ , if  $[y_{\mathcal{O}_\eta}^{M,\text{new}}]$  is a new optimal solution of Eq. (2.24), the global accuracy of the identification stems from the comparison between  $\mathcal{C}([ \eta^{\text{exp,new}}, [y_{\mathcal{O}_\eta}^M], [\Psi^*] )$  and  $\mathcal{C}([ \eta^{\text{exp,new}}, [y_{\mathcal{O}_\eta}^{M,\text{new}}], [\Psi^*] )$ .

**2.4. Identification of the PCE truncation parameters.** As shown in Introduction, two truncation parameters,  $N_g$  and  $p$ , appear in the truncated PCE,  $\boldsymbol{\eta}^{\text{chaos}}(N) = [y]\boldsymbol{\Psi}(\boldsymbol{\xi}, p)$ , of  $\boldsymbol{\eta}$ . The values of these parameters have to be determined from a convergence analysis. The objective of this section is thus to give the fundamental elements to perform such a convergence analysis.

**2.4.1. Definition of a log error function.** For each component  $\eta_k^{\text{chaos}}(N)$  of the truncated PCE,  $\boldsymbol{\eta}^{\text{chaos}}(N) = [y]\boldsymbol{\Psi}(\boldsymbol{\xi}, p)$ , of  $\boldsymbol{\eta}$ , the  $L^1$ -log error function  $err_k$  is introduced as described in [29]:

$$\forall 1 \leq k \leq N_\eta, \quad err_k(N_g, p) = \int_{\text{BI}_k} |\log_{10}(p_{\eta_k}(x_k)) - \log_{10}(p_{\eta_k^{\text{chaos}}}(x_k))| dx_k, \quad (2.25)$$

where:

- $\text{BI}_k$  is the support of  $\eta_k^{\text{exp}}$ ;
- $p_{\eta_k}$  and  $p_{\eta_k^{\text{chaos}}}$  are the PDF of  $\eta_k$  and  $\eta_k^{\text{chaos}}$  respectively.

The multidimensional error function  $err(N_g, p)$  is then deduced from the unidimensional  $L^1$ -log error function as:

$$err(N_g, p) = \sum_{k=1}^{N_\eta} err_k(N_g, p). \quad (2.26)$$

The parameters  $N_g$  and  $p$  have thus to be determined to minimize the multidimensional  $L^1$ -log error function  $err(N_g, p)$ .

For given values of truncation parameters  $N_g$  and  $p$ , it is reminded that PCE coefficients matrix  $[y]$  is searched in order to maximize the multidimensional log-likelihood function, which allows us to consider a priori strongly correlated problems. Once this matrix  $[y]$  is identified, it is possible to generate as many independent realizations of truncated PCE  $\boldsymbol{\eta}^{\text{chaos}}(N)$  as needed to estimate as precisely as possible the non parametric estimator  $\hat{p}_{\mathcal{U}}$  of its multidimensional PDF. The number  $\nu^{\text{exp}}$  of available experimental realizations of  $\boldsymbol{\eta}$  is however limited. This number is generally too small for the non parametric estimator of multidimensional PDF  $p_{\boldsymbol{\eta}}$  of  $\boldsymbol{\eta}$  to be relevant, whereas it is most of the time large enough to define the estimators of the marginals of  $p_{\boldsymbol{\eta}}$ . Therefore, the log-error functions defined by Eqs. (2.25) and (2.26)

only consider the marginals of the PDF of  $p_\eta$  and  $p_\eta^{\text{chaos}}$ . In addition, the logarithm function has been introduced in order to measure the errors of the very small values of the probability density function (the tails of the probability density function).

#### 2.4.2. Definition of an admissible set for the truncation parameters.

As it exists an isoprobabilistic transformation between  $\eta$  and  $(\Xi_1, \dots, \Xi_{N_\eta})$ , where  $\{\Xi_k, 1 \leq k \leq N_\eta\}$  is a set of  $N_\eta$  independent centered normalized Gaussian random variables, the convergence analysis can be restricted to the values of  $N_g$  which verify:

$$N_g \leq N_\eta. \quad (2.27)$$

Moreover, imposing the  $(N_\eta \times N)$  real matrix  $[y]$  to be in  $\mathcal{O}_\eta$  amounts to imposing  $\frac{N_\eta(N_\eta+3)}{2}$  constraints on  $[y]$ , which implies:

$$N_\eta N \geq \frac{N_\eta(N_\eta+3)}{2} \Leftrightarrow N \geq \frac{N_\eta+3}{2}. \quad (2.28)$$

However, the algorithms developed in [29], on which the solving of the optimization problem, defined by Eq. (2.24), is based, need the more restrictive condition:

$$N \geq N_\eta + 1. \quad (2.29)$$

We will therefore consider  $\mathcal{Q}(N_\eta)$  the set of the admissible values for  $p$  and  $N_g$  with:

$$\mathcal{Q}(N_\eta) = \{(p, N_g) \in \mathbb{N}^2, \mid N_g \leq N_\eta, N = (N_g + p)! / (N_g! p!) \geq N_\eta + 1\}. \quad (2.30)$$

Theoretically, increasing  $p$  and  $N_g$  adds terms in the PCE of the considered random vector, and therefore should induce the decrease of the error function:

$$\forall p^* \geq p, N_g^* \geq N_g, \text{err}(N_g, p) \geq \max\{\text{err}(N_g^*, p), \text{err}(N_g, p^*)\} \geq \min\{\text{err}(N_g^*, p), \text{err}(N_g, p^*)\} \geq \text{err}(N_g^*, p^*). \quad (2.31)$$

However, the higher the values of  $p$  and  $N_g$  are, the bigger the PCE coefficients matrix is, the harder the numerical identification is. Hence, introducing  $\varepsilon$  as an error threshold, which has to be adapted to the problem, let  $\mathcal{P}(\varepsilon, N_\eta)$  be the set:

$$\mathcal{P}(\varepsilon, N_\eta) = \{(p, N_g) \in \mathcal{Q}(N_\eta) \mid \text{err}(N_g, p) \leq \varepsilon\}. \quad (2.32)$$

Finally, given the error threshold  $\varepsilon$ , rather than directly minimizing the  $L^1$ -log error function  $\text{err}(N_g, p)$ , it appears to be more accurate to look for the optimal values of  $p$  and  $N_g$  that minimize the size of the projection basis  $N = (N_g + p)! / (N_g! p!)$ :

$$(p, N_g) = \arg \min_{(p^*, N_g^*) \in \mathcal{P}(\varepsilon, N_\eta)} (N_g^* + p^*)! / (N_g^*! p^*!). \quad (2.33)$$

If the polynomial order (which is a priori unknown) of the non truncated PCE of  $\eta$  is infinite, it may not exist values of  $p$  and  $N_g$  in  $\mathcal{P}(\varepsilon, N_\eta)$  for error function  $\text{err}(N_g, p)$  to be inferior to small values of  $\varepsilon$ . In this case, the former algorithms can nevertheless be used to find the most accurate values of  $p$  and  $N_g$  with respect to an available computational cost.

**3. Adaptation of the PCE identification method in high dimension.** As it has been presented in the former sections, the  $(N \times \nu^{\text{chaos}})$  real matrix  $[\Psi]$  gathers  $\nu^{\text{chaos}}$  independent realizations of the  $N$ -dimension PCE basis  $\{\psi_{\alpha}(\xi), \alpha \in \mathcal{A}_p\}$ . Eqs. (2.23) and (2.24) underline the fact that the numerical identification of the PCE coefficients  $[y]$  can be seen as the minimization of a cost function involving the elements of the  $(N_{\eta} \times \nu^{\text{chaos}})$  real matrix of independent realizations  $[U] = [y][\Psi]$  and the elements of the  $(N_{\eta} \times \nu^{\text{exp}})$  real matrix  $[\eta^{\text{exp}}] = [\boldsymbol{\eta}^{(1)} \dots \boldsymbol{\eta}^{(\nu^{\text{exp}})}]$ . In theoretical terms, this cost function should be minimum when the multidimensional PDF  $p_U$  of  $\mathbf{U} = [y]\Psi(\xi, p)$  is as near as possible to the multidimensional PDF  $p_{\boldsymbol{\eta}}$  of  $\boldsymbol{\eta}$ . In practical terms, this cost function is however minimum when  $\hat{p}_U$  is as near as possible to  $\hat{p}_{\boldsymbol{\eta}}$ , where  $\hat{p}_U$  and  $\hat{p}_{\boldsymbol{\eta}}$  are the multidimensional non parametric estimators of  $p_U$  and  $p_{\boldsymbol{\eta}}$  defined by Eq. (2.18). With respect to  $\nu^{\text{exp}}$  and  $\nu^{\text{chaos}}$ , three bias are then introduced in the PCE identification:

- a bias due to a lack of information on  $\boldsymbol{\eta}$ :

$$b^{(1)}(\nu^{\text{exp}}) = \int_{\mathbb{R}^{N_{\eta}}} |\hat{p}_{\boldsymbol{\eta}}(\mathbf{x}) - p_{\boldsymbol{\eta}}(\mathbf{x})| d\mathbf{x}, \quad (3.1)$$

- a bias due to a lack of information on  $\mathbf{U}$ :

$$b^{(2)}(\nu^{\text{chaos}}) = \int_{\mathbb{R}^{N_{\eta}}} |\hat{p}_U(\mathbf{x}) - p_U(\mathbf{x})| d\mathbf{x}, \quad (3.2)$$

- a bias due to the truncation and to the fact that the global maximum is not necessary reached:

$$b^{(3)}(\nu^{\text{exp}}, \nu^{\text{chaos}}) = \int_{\mathbb{R}^{N_{\eta}}} |\hat{p}_{\boldsymbol{\eta}}(\mathbf{x}) - \hat{p}_U(\mathbf{x})| d\mathbf{x}. \quad (3.3)$$

These three bias could also be expressed with respect to the statistical moments of  $\boldsymbol{\eta}$  and  $\mathbf{U}$ . For instance, when focusing on the autocorrelation matrix, let  $err^1$ ,  $err^2$  and  $err^3$  be the autocorrelation errors corresponding respectively to the bias  $b^{(1)}$ ,  $b^{(2)}$  and  $b^{(3)}$ :

$$err^1(\nu^{\text{exp}}) = \left\| [R_{\boldsymbol{\eta}}] - [\hat{R}_{\boldsymbol{\eta}}(\nu^{\text{exp}})] \right\|_F / \|[R_{\boldsymbol{\eta}}]\|_F, \quad (3.4)$$

$$err^2(\nu^{\text{chaos}}) = \left\| [R_{\boldsymbol{\eta}}^{\text{chaos}}(N)] - [\hat{R}_U(\nu^{\text{chaos}})] \right\|_F / \|[R_{\boldsymbol{\eta}}^{\text{chaos}}(N)]\|_F, \quad (3.5)$$

$$err^3(\nu^{\text{exp}}, \nu^{\text{chaos}}) = \left\| [\hat{R}_U(\nu^{\text{chaos}})] - [\hat{R}_{\boldsymbol{\eta}}(\nu^{\text{exp}})] \right\|_F / \left\| [\hat{R}_{\boldsymbol{\eta}}(\nu^{\text{exp}})] \right\|_F, \quad (3.6)$$

where  $\|\cdot\|_F$  is the Frobenius norm of matrices, and where it is reminded from Eqs. (2.8) and (2.13) that:

$$\begin{cases} [\hat{R}_{\boldsymbol{\eta}}(\nu^{\text{exp}})] = \frac{1}{\nu^{\text{exp}}} [\eta^{\text{exp}}][\eta^{\text{exp}}]^T, \\ [\hat{R}_U(\nu^{\text{chaos}})] = \frac{1}{\nu^{\text{chaos}}} [U][U]^T = [y] \left( \frac{1}{\nu^{\text{chaos}}} [\Psi][\Psi]^T \right) [y]^T, \\ [R_{\boldsymbol{\eta}}^{\text{chaos}}(N)] = [y][y]^T. \end{cases} \quad (3.7)$$

Hence, the smaller these three errors are, the more precise the PCE identification is. The  $\nu^{\text{exp}}$  independent realizations  $\{\boldsymbol{\eta}^{(1)}, \dots, \boldsymbol{\eta}^{(\nu^{\text{exp}})}\}$  being the maximum available information on  $\boldsymbol{\eta}$ , the bias  $b^{(1)}$  and the autocorrelation error  $err^1$  cannot be decreased, whereas the set  $\mathcal{O}_\eta$ , which was introduced to guarantee that  $[R_\eta^{\text{chaos}}(N)] = [\hat{R}_\eta(\nu^{\text{exp}})]$ , aims at reducing  $b^{(2)}$ ,  $b^{(3)}$ ,  $err^2$  and  $err^3$ . Therefore, imposing  $[y]$  to be in  $\mathcal{O}_\eta$  leads us to:

$$\begin{aligned} err^2(\nu^{\text{chaos}}) &= err^3(\nu^{\text{exp}}, \nu^{\text{chaos}}) \\ &= \left\| [y] \left( \frac{1}{\nu^{\text{chaos}}} [\Psi][\Psi]^T - [I_N] \right) [y]^T \right\|_F / \left\| [y][I_N][y]^T \right\|_F. \end{aligned} \quad (3.8)$$

The following asymptotical property can thus be deduced from Eq. (1.12):

$$\lim_{\nu^{\text{chaos}} \rightarrow +\infty} err^2(\nu^{\text{chaos}}) = \lim_{\nu^{\text{chaos}} \rightarrow +\infty} err^3(\nu^{\text{exp}}, \nu^{\text{chaos}}) = 0, \quad (3.9)$$

which is equivalent to say that the larger  $\nu^{\text{chaos}}$  is, the more accurate the PCE identification should be. However, from a practical point of view, the value of  $\nu^{\text{chaos}}$  is fixed by the available computation resources. As an extension of the work presented in [30], this section aims at quantifying the divergence of the ratio:

$$r = \left\| \frac{1}{\nu^{\text{chaos}}} [\Psi][\Psi]^T - [I_N] \right\|_F / \left\| [I_N] \right\|_F, \quad (3.10)$$

when the truncation parameters  $N_g$  and  $p$  increase for several statistical measures. From Eq. (3.8),  $r$ , defined by Eq. (3.10), can be seen as a general characterization of the autocorrelation errors  $err^2$  and  $err^3$ . This divergence being very detrimental to the PCE identification in high dimension, a new decomposition of the PCE coefficient matrix  $[y]$  will be then presented in this section to make  $err^2$  and  $err^3$  be zero for any value of  $N_g$  and  $p$ .

**3.1. Decomposition of the matrix of independent realizations.** To better emphasize the influence of the truncation parameters on the ratio  $r$ , a rewriting of the matrix  $[\Psi]$  is first presented.

**3.1.1. Theoretical basis of the decomposition.** From Eq. (1.10), matrix  $[\Psi]$  gathers  $\nu^{\text{chaos}}$  columns  $\{\Psi(\boldsymbol{\xi}(\theta_n), p), 1 \leq n \leq \nu^{\text{chaos}}\}$ , which are independent realizations of the  $N$ -dimension PCE basis  $\{\psi_\alpha(\boldsymbol{\xi}), \alpha \in \mathcal{A}_p\}$ . This basis being orthonormal leads us to the asymptotical condition on  $[\Psi]$ , defined by Eq. (1.12). Moreover, Eq. (1.6) implies that  $[\Psi]$  can be expressed as:

$$[\Psi] = [A][M], \quad (3.11)$$

where  $[A]$  is the  $(N \times N)$  real matrix that gathers the coefficients of the orthonormal polynomials with respect to the probability measure of the  $N_g$ -dimension PCE germ,  $\boldsymbol{\xi} = (\xi_1, \dots, \xi_{N_g})$ , and  $[M]$  is a  $(N \times \nu^{\text{chaos}})$  real matrix of  $\nu^{\text{chaos}}$  independent realizations of the multi-index monomials  $\mathcal{M}_\alpha(\boldsymbol{\xi}) = \xi_1^{\alpha_1} \times \dots \times \xi_{N_g}^{\alpha_{N_g}}$ , for any value  $\alpha$  in  $\mathcal{A}_p$ :

$$[M] = [\mathcal{E}(\xi(\theta_1), p) \quad \cdots \quad \mathcal{E}(\xi(\theta_{\nu^{\text{chaos}}}), p)], \quad (3.12)$$

$$\mathcal{E}(\xi, p) = (\mathcal{M}_{\alpha^{(1)}}(\xi), \dots, \mathcal{M}_{\alpha^{(N)}}(\xi)). \quad (3.13)$$

If  $[A]$  is independent of  $[M]$ , Eq. (3.11) certifies that, if the columns of  $[M]$  are independent, then the columns of  $[\Psi]$  stay independent. Let  $[R_{\mathcal{E}}]$  be the autocorrelation matrix of the random vector  $\mathcal{E}(\xi, p)$ :

$$[R_{\mathcal{E}}] = E \left( \mathcal{E}(\xi, p) \mathcal{E}(\xi, p)^T \right). \quad (3.14)$$

It can be deduced from Eqs. (1.12), (3.11), (3.12) and (3.14) that:

$$[R_{\mathcal{E}}] = \lim_{\nu^{\text{chaos}} \rightarrow +\infty} \frac{1}{\nu^{\text{chaos}}} [M][M]^T = [A]^{-1} [A]^{-T}. \quad (3.15)$$

According to this decomposition, computing the classical Gram-Schmidt orthogonalization to identify the polynomial basis coefficients only requires the calculation of  $[A]^{-T}$ , which corresponds to the Cholesky decomposition matrix of the positive definite matrix  $[R_{\mathcal{E}}]$ . Hence, by construction, the matrix  $[\Psi]$  can be written as the product of a lower triangular matrix  $[A]$  and a matrix  $[M]$  of independent realizations of a multi-index random vector  $\mathcal{E}(\xi, p)$ .

**3.1.2. Practical computation of matrix  $[\Psi]$ .** Thanks to Eq. (3.11), matrix  $[\Psi]$  can be numerically computed without requiring computational recurrence formula nor algebraic explicit representation. An illustration of the method is presented hereinafter for a PCE based on a Gaussian measure. This development can be directly extended to any value of  $p$  and  $N_g$ , as well as to other statistical measures. Let  $\xi_1$  and  $\xi_2$  be two independent normalized Gaussian random variables, such that  $\xi = (\xi_1, \xi_2)$ , and  $\alpha = (\alpha_1, \alpha_2)$ . Choosing  $p = 2$  and  $N_g = 2$ , which corresponds to  $N = 6$ , leads to the following definition of  $\mathcal{E}(\xi, p)$ :

$$\mathcal{E}(\xi, 2) = (1, \xi_1, \xi_2, \xi_1 \xi_2, \xi_1^2, \xi_2^2). \quad (3.16)$$

According to this equation, matrix  $[M]$  can thus be easily deduced from  $\nu^{\text{chaos}}$  independent realizations of  $\xi$ . Moreover, let  $[\alpha]$  be the  $(N_g \times N)$  real matrix which gathers the admissible values for  $\alpha$  in  $\mathcal{A}_p$ :

$$[\alpha] = \begin{bmatrix} 0 & 1 & 0 & 1 & 2 & 0 \\ 0 & 0 & 1 & 1 & 0 & 2 \end{bmatrix} \leftrightarrow \mathcal{A}_p = \{(0, 0), (1, 0), (0, 1), (1, 1), (2, 0), (0, 2)\}. \quad (3.17)$$

The random variables  $\xi_1$  and  $\xi_2$  being independent, normalized and Gaussian, the autocorrelation matrix  $[R_{\mathcal{E}}]$  can thus be written as:

$$\begin{aligned} \forall i, j \in \{1, \dots, N\}, [R_{\mathcal{E}}]_{ij} &= E \left( \xi_1^{[\alpha]_{1i} + [\alpha]_{1j}} \times \dots \times \xi_{N_g}^{[\alpha]_{N_g i} + [\alpha]_{N_g j}} \right) \\ &= E \left( \xi_1^{[\alpha]_{1i} + [\alpha]_{1j}} \right) \times \dots \times E \left( \xi_{N_g}^{[\alpha]_{N_g i} + [\alpha]_{N_g j}} \right), \end{aligned} \quad (3.18)$$

where, for  $1 \leq \ell \leq N_g$ :

$$\begin{cases} E(\xi_\ell^q) = 0 & \text{if } q \text{ is not even,} \\ E(\xi_\ell^q) = \frac{q!}{(q/2)!2^{q/2}} & \text{if } q \text{ is even.} \end{cases} \quad (3.19)$$

Therefore, Eq. (3.15) allows us to numerically find back in  $[A]$  the multidimensional Hermite polynomials  $H_{\alpha_1} \times \dots \times H_{\alpha_{N_g}}$ :

$$\forall x \in \mathbb{R}, \begin{cases} H_0(x_1) \times H_0(x_2) = 1 \\ H_1(x_1) \times H_0(x_2) = x_1 \\ H_0(x_1) \times H_1(x_2) = x_2 \\ H_1(x_1) \times H_1(x_2) = x_1 x_2 \\ H_2(x_1) \times H_0(x_2) = \frac{x_1^2 - 1}{\sqrt{2}} \\ H_0(x_1) \times H_2(x_2) = \frac{x_2^2 - 1}{\sqrt{2}} \end{cases} \leftrightarrow [A] = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ \frac{-1}{\sqrt{2}} & 0 & 0 & 0 & \frac{1}{\sqrt{2}} & 0 \\ \frac{-1}{\sqrt{2}} & 0 & 0 & 0 & 0 & \frac{1}{\sqrt{2}} \end{pmatrix}. \quad (3.20)$$

Noticing that:

- if  $\xi$  is a random variable uniformly distributed on  $[-1, 1]$ :

$$\begin{cases} E(\xi^q) = 0 & \text{if } q \text{ is not even,} \\ E(\xi^q) = \frac{1}{q+1} & \text{if } q \text{ is even,} \end{cases} \quad (3.21)$$

- if the random variable  $\xi$  is a random variable that is characterized by a normalized exponential distribution on  $[0, +\infty[$ :

$$E(\xi^q) = q!, \quad (3.22)$$

this method can directly be generalized to the uniform and exponential cases to compute the multidimensional Legendre and Laguerre polynomial coefficients, but also to an arbitrary probability measure for the germ  $\xi$ .

**3.2. Influence of the truncation parameters and of the choice for the PCE probability measure.** The convergence properties of ratio  $r$  when  $\nu^{\text{chaos}}$  tends to infinity are strongly related to the statistical properties of germ  $\xi$ . This section aims therefore to emphasize the dominant trends of this specific link, and to highlight the difficulties brought about by the divergence of ratio  $r$ , when trying to perform analysis of convergence in high dimension.

The definition of the Frobenius norm allows us to write that:

$$r = \left\| \frac{1}{\nu^{\text{chaos}}} [\Psi][\Psi]^T - [I_N] \right\|_F / \|[I_N]\|_F = \sqrt{N} \Sigma(\nu^{\text{chaos}}), \quad (3.23)$$

where  $\Sigma(\nu^{\text{chaos}})$  is such that:

$$\{\Sigma(\nu^{\text{chaos}})\}^2 = \frac{1}{N^2} \sum_{1 \leq i, j \leq N} \left( \left( \frac{1}{\nu^{\text{chaos}}} [\Psi][\Psi]^T - [I_N] \right)_{ij} \right)^2. \quad (3.24)$$

By construction,  $\{\Sigma(\nu^{\text{chaos}})\}^2$  is an assessment of the mean value of the squared difference between the elements of  $\frac{1}{\nu^{\text{chaos}}}[\Psi][\Psi]^T$  and the elements of the identity matrix  $[I_N]$ . Hence, if  $\{\Sigma(\nu^{\text{chaos}})\}^2$  remains constant when the size  $N$  of the polynomial basis increases, the ratio  $r$  should increase as  $\sqrt{N}$ . Moreover, Eqs. (3.11) and (3.15) yield,

$$\frac{1}{\nu^{\text{chaos}}}[\Psi][\Psi]^T - [I_N] = [A] \left( \frac{1}{\nu^{\text{chaos}}}[M][M]^T - [R_{\mathcal{E}}] \right) [A]^T. \quad (3.25)$$

For all  $(i, j)$  in  $\{1, \dots, N\}^2$ ,  $[R_{\mathcal{E}}]_{ij}$  is such that:

$$[R_{\mathcal{E}}]_{ij} = E \left( \xi_1^{\alpha_1^{(i)} + \alpha_1^{(j)}} \times \dots \times \xi_{N_g}^{\alpha_{N_g}^{(i)} + \alpha_{N_g}^{(j)}} \right). \quad (3.26)$$

Let  $[\widehat{R}_{\mathcal{E}}]$  be the following estimator of  $[R_{\mathcal{E}}]$ :

$$[\widehat{R}_{\mathcal{E}}]_{ij} = \frac{1}{\nu^{\text{chaos}}} \sum_{n=1}^{\nu^{\text{chaos}}} \left( \Xi_1^{(n)} \right)^{\alpha_1^{(i)} + \alpha_1^{(j)}} \times \dots \times \left( \Xi_{N_g}^{(n)} \right)^{\alpha_{N_g}^{(i)} + \alpha_{N_g}^{(j)}}, \quad (3.27)$$

where  $\{\Xi^{(n)} = (\Xi_1^{(n)}, \dots, \Xi_{N_g}^{(n)}), 1 \leq n \leq \nu^{\text{chaos}}\}$  is a set of  $\nu^{\text{chaos}}$  independent  $N_g$ -dimension random vectors, which have the same PDF than  $\xi$ . The central limit theorem yields that, for all  $(i, j)$  in  $\{1, \dots, N\}^2$ , we have:

$$\sqrt{\frac{\nu^{\text{chaos}}}{\text{Var} \left( \xi_1^{\alpha_1^{(i)} + \alpha_1^{(j)}} \times \dots \times \xi_{N_g}^{\alpha_{N_g}^{(i)} + \alpha_{N_g}^{(j)}} \right)}} \left( [\widehat{R}_{\mathcal{E}}]_{ij} - [R_{\mathcal{E}}]_{ij} \right) \xrightarrow{\text{law}} \mathcal{N}(0, 1), \quad (3.28)$$

where  $\mathcal{N}(0, 1)$  is the normalized Gaussian distribution, and  $\text{Var}(\cdot)$  is the variance. Under this formalism, it can be noticed that  $\frac{1}{\nu^{\text{chaos}}}[\Psi][\Psi]^T$  is one particular realization of  $[\widehat{R}_{\mathcal{E}}]$ . Hence, from Eqs. (3.24), (3.25) and (3.28), we deduce that:

- if  $\text{Var} \left( \xi_1^{\alpha_1^{(i)} + \alpha_1^{(j)}} \times \dots \times \xi_{N_g}^{\alpha_{N_g}^{(i)} + \alpha_{N_g}^{(j)}} \right) \leq \text{Var} \left( \xi_1^{\alpha_1^{(i)} + \alpha_1^{(j)}} \times \dots \times \xi_{N_g+1}^{\alpha_{N_g+1}^{(i)} + \alpha_{N_g+1}^{(j)}} \right)$ , then  $\Sigma(\nu^{\text{chaos}})$  potentially increases with respect to  $N_g$ .
- if  $\text{Var} \left( \xi_{\ell}^{\alpha_{\ell}^{(i)}} \right) \leq \text{Var} \left( \xi_{\ell}^{\alpha_{\ell}^{(j)}} \right)$  for  $\alpha_{\ell}^{(i)} \leq \alpha_{\ell}^{(j)}$ , then  $\Sigma(\nu^{\text{chaos}})$  potentially increases with respect to  $p$ .

As an illustration, for each couple  $(N_g, p)$  such that  $1 \leq p \leq 10$  and  $1 \leq N_g \leq 6$ , three sets,  $\{[\Psi_U^{(m)}(p, N_g)], 1 \leq m \leq 1000\}$ ,  $\{[\Psi_G^{(m)}(p, N_g)], 1 \leq m \leq 1000\}$  and  $\{[\Psi_E^{(m)}(p, N_g)], 1 \leq m \leq 1000\}$ , are computed, such that  $[\Psi_U^{(m)}(p, N_g)]$ ,  $[\Psi_G^{(m)}(p, N_g)]$  and  $[\Psi_E^{(m)}(p, N_g)]$  refer to particular  $(N \times \nu^{\text{chaos}})$  real matrices of independent realisations of the basis  $\{\psi_{\alpha}, \alpha(\xi) \in \mathcal{A}_p\}$ , in the uniform, the Gaussian and the exponential cases, respectively. Hence, defining:

$$\begin{cases} r_U^m(\nu^{\text{chaos}}) = \left\| \frac{1}{\nu^{\text{chaos}}} [\Psi_U^{(m)}(p, N_g)] [\Psi_U^{(m)}(p, N_g)]^T - [I_N] \right\|_F / \|[I_N]\|_F, \\ r_G^m(\nu^{\text{chaos}}) = \left\| \frac{1}{\nu^{\text{chaos}}} [\Psi_G^{(m)}(p, N_g)] [\Psi_G^{(m)}(p, N_g)]^T - [I_N] \right\|_F / \|[I_N]\|_F, \\ r_E^m(\nu^{\text{chaos}}) = \left\| \frac{1}{\nu^{\text{chaos}}} [\Psi_E^{(m)}(p, N_g)] [\Psi_E^{(m)}(p, N_g)]^T - [I_N] \right\|_F / \|[I_N]\|_F, \end{cases} \quad (3.29)$$

allows us to compute, in each case, three approximations  $err_U^{\text{ortho}}(p, N_g)$ ,  $err_G^{\text{ortho}}(p, N_g)$  and  $err_E^{\text{ortho}}(p, N_g)$  of the mean value of the ratio  $r$ , defined in Eq. (3.23), such that:

$$\begin{cases} err_U^{\text{ortho}}(p, N_g) = \frac{1}{1000} \sum_{m=1}^{1000} r_U^m(\nu^{\text{chaos}}), \\ err_G^{\text{ortho}}(p, N_g) = \frac{1}{1000} \sum_{m=1}^{1000} r_G^m(\nu^{\text{chaos}}), \\ err_E^{\text{ortho}}(p, N_g) = \frac{1}{1000} \sum_{m=1}^{1000} r_E^m(\nu^{\text{chaos}}). \end{cases} \quad (3.30)$$

For  $\nu^{\text{chaos}} = 1000$ , in figure 3.1, the two factors which make the ratio  $r$  diverge with respect to  $p$  and  $N_g$  can therefore be emphasized. On the first hand, if increasing  $p$  or  $N_g$  does not increase the variance of the elements of  $\mathcal{E}(\xi, p)$ , which is the case if the PCE germ  $\xi$  is characterized by an uniform distribution (see Eq. (3.21)), the ratio  $r$  increases approximately as  $\sqrt{N}$ . On the other hand, if increasing  $p$  or  $N_g$  increases the variance of the element of  $\mathcal{E}(\xi, p)$ , as it is the case if the PCE germ  $\xi$  is characterized by a Gaussian or exponential distribution (see Eqs. (3.19) and (3.22)), the ratio  $r$  diverges very quickly with respect to the truncation parameters, and bias the PCE identification results.

As a conclusion, for a fixed value of  $\nu^{\text{chaos}}$ , the difference  $\frac{1}{\nu^{\text{chaos}}} [\Psi][\Psi]^T - [I_N]$  increases when  $p$  and  $N_g$  increase. Therefore, imposing  $[y]$  to be in  $\mathcal{O}_\eta$  introduces a numerical bias in the PCE identification, which becomes very important when high values of  $p$  and  $N_g$  are needed. Such a phenomenon prevents thus to perform the analysis of convergence of the PCE in high dimension, especially when dealing with Gaussian and exponential PCE germs.

**3.3. Adaptation of the optimization problem.** In this section, fixed values for  $\nu^{\text{chaos}}$ ,  $p$  and  $N_g$  are considered. According to the notations of Section 3.1, a  $(N \times \nu^{\text{chaos}})$  real matrix of independent realizations  $[\Psi] = [A][M]$  can then be constructed. Under the condition  $\nu^{\text{chaos}} \geq N$ ,  $\frac{1}{\nu^{\text{chaos}}} [M][M]^T$  is positive definite by construction, which allows writing:

$$\frac{1}{\nu^{\text{chaos}}} [M][M]^T = [L][L]^T, \quad (3.31)$$

where  $[L]$  is the Cholesky decomposition of  $\frac{1}{\nu^{\text{chaos}}} [M][M]^T$ , which yields:

$$\frac{1}{\nu^{\text{chaos}}} [\Psi][\Psi]^T = [A][L][L]^T[A]^T = [B][B]^T, \quad (3.32)$$

$$[B] = [A][L]. \quad (3.33)$$

The matrix:



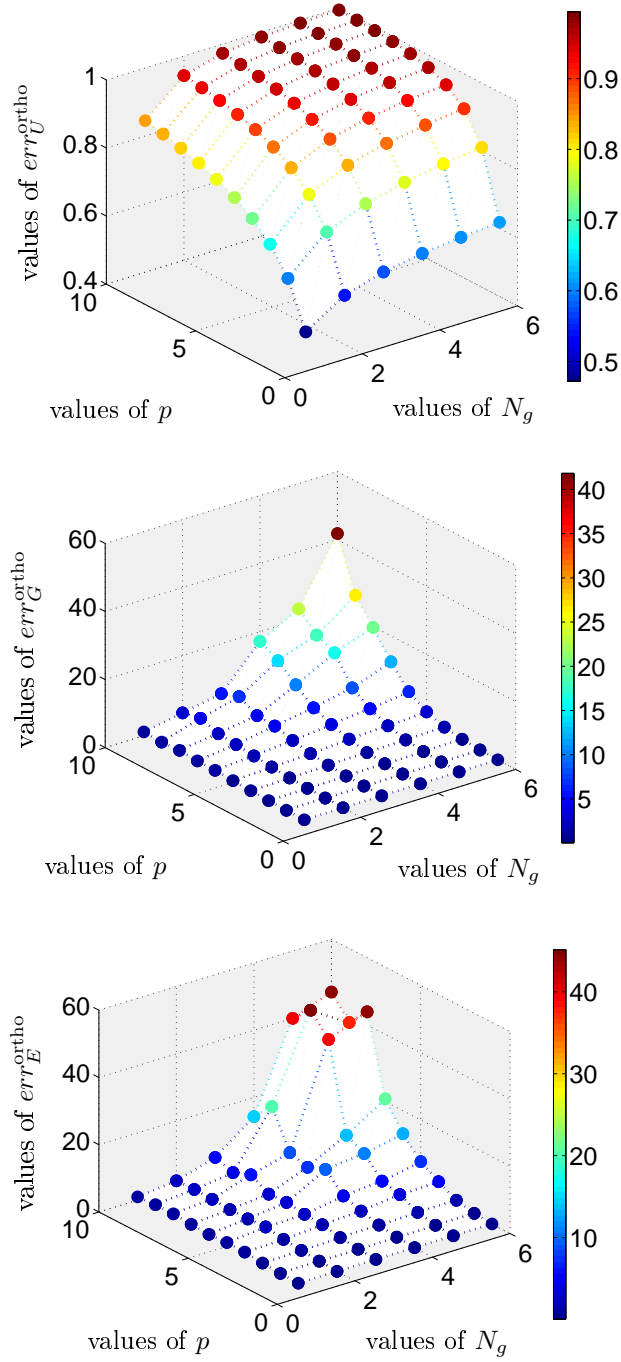


Figure 3.1: Graphs of the errors  $err_U^{\text{ortho}}$ ,  $err_G^{\text{ortho}}$  and  $err_E^{\text{ortho}}$  with respect to the truncation parameters  $N_g$  and  $p$ .

$$[\tilde{\Psi}] = [B]^{-1}[\Psi], \quad (3.34)$$

is then introduced, such that by construction:

$$\frac{1}{\nu^{\text{chaos}}} [\tilde{\Psi}] [\tilde{\Psi}]^T = [I_N]. \quad (3.35)$$

Using the notations of Section 2, let  $[y^*]$  be a  $(N_\eta \times N)$  real matrix such that the random vector  $\mathbf{U}$  is defined as:

$$\mathbf{U} = [y^*] \Psi(\xi, p). \quad (3.36)$$

Hence,  $\nu^{\text{chaos}}$  independent realizations of  $\mathbf{U}$  can be directly deduced from the matrix  $[\Psi]$  and gathered in the matrix  $[U] = [y^*][\Psi]$ . Defining  $[z]$  such that:

$$[z] = [y^*][B], \quad (3.37)$$

therefore yields the equality:

$$[U] = [y^*][\Psi] = ([z][B]^{-1}) ([B][\tilde{\Psi}]) = [z][\tilde{\Psi}]. \quad (3.38)$$

If  $[z]$  is in  $\mathcal{O}_\eta$ ,  $[z][z]^T = [\hat{R}_\eta(\nu^{\text{exp}})]$ , which implies that:

$$\begin{cases} [\hat{R}_U(\nu^{\text{chaos}})] = \frac{1}{\nu^{\text{chaos}}} [U][U]^T = [z] \left( \frac{1}{\nu^{\text{chaos}}} [\tilde{\Psi}][\tilde{\Psi}]^T \right) [z]^T = [z][z]^T = [\hat{R}_\eta(\nu^{\text{exp}})], \\ [R_U] = E(\mathbf{U}\mathbf{U}^T) = \lim_{\nu^{\text{chaos}} \rightarrow \infty} [\hat{R}_U(\nu^{\text{chaos}})] = [\hat{R}_\eta(\nu^{\text{exp}})]. \end{cases} \quad (3.39)$$

From Eqs. (3.5) and (3.6), it can thus be deduced that imposing  $[z]$  to be an element of  $\mathcal{O}_\eta$  guarantees that, for any  $\nu^{\text{chaos}} \geq N$ , we have  $\text{err}^2(\nu^{\text{chaos}}) = \text{err}^3(\nu^{\text{exp}}, \nu^{\text{chaos}}) = 0$ .

Hence, whereas the optimization problem defined by Eq. (2.24) is perturbed by autocorrelation errors, the new optimization problem:

$$\begin{cases} \begin{bmatrix} y_{\mathcal{O}_\eta}^M \end{bmatrix} = \begin{bmatrix} z_{\mathcal{O}_\eta}^M \end{bmatrix} [B^{-1}], \\ \begin{bmatrix} z_{\mathcal{O}_\eta}^M \end{bmatrix} = \arg \max_{[z^*] \in \mathcal{W}} \mathcal{C} \left( [\eta^{\text{exp}}], [z^*], [\tilde{\Psi}] \right), \end{cases} \quad (3.40)$$

is no more affected, which allows us to consider high values of the truncation parameters  $N_g$  and  $p$ . Equation (3.38) underlines that the two former optimization problems are equivalent, as the independent realizations of  $\mathbf{U}$  have just been rewritten. Only the research set, for the PCE coefficient matrix, has been modified, which allows the numerical bias due to the finite dimension of  $[\Psi]$  to be reduced.

Finally, if  $[y]$  is the coefficients matrix of the truncated PCE,  $\boldsymbol{\eta}^{\text{chaos}}(N)$ , of random vector  $\boldsymbol{\eta}$ , such that  $\boldsymbol{\eta}^{\text{chaos}}(N) = [y] \Psi(\xi, p)$ , a good estimation of  $[y]$  in high dimension can be computed by solving the optimization problem defined by Eq. (3.40).

**3.4. Remarks on the new optimization problem.** It has to be noticed that  $[\tilde{\Psi}]$  is unique, and keeps exactly the same structure than  $[\Psi]$ . Indeed, let  $[L^{\text{asym}}] = [A]^{-1}$  be the Cholesky decomposition matrix of the autocorrelation matrix  $[R_{\mathcal{E}}]$ , which is defined by Eq. (3.14). Hence, from Eq. (3.11),  $[\Psi] = [L^{\text{asym}}]^{-1}[M]$ , which has to be compared to  $[\tilde{\Psi}] = [B]^{-1}[\Psi] = ([L]^{-1}[A]^{-1})([A][M]) = [L]^{-1}[M]$ , where  $[L]$  and  $[L^{\text{asym}}]$  are two lower triangular matrices. Whereas  $[L^{\text{asym}}]$  implies the asymptotical orthonormality,  $[L]$  guarantees the numerical orthonormality. Moreover, from Eq. (3.40), the optimal PCE coefficients matrix  $[y]$  is approximated as a product of two matrices:

$$[y] \approx [z_{\mathcal{O}_\eta}^M] [B]^{-1}. \quad (3.41)$$

For a fixed value of  $N$ ,  $[B]$  is strongly dependent on  $\nu^{\text{chaos}}$  and  $[\Psi]$ . From Eq. (1.12), it also verifies the asymptotical property:

$$\lim_{\nu^{\text{chaos}} \rightarrow \infty} [B] = [I_N], \quad (3.42)$$

which implies that  $[z_{\mathcal{O}_\eta}^M]$  converges towards  $[y]$  if sufficiently high values of  $\nu^{\text{chaos}}$  is considered. Hence, the less dependent on  $[\Psi]$  the matrix  $[z_{\mathcal{O}_\eta}^M]$  is, the more accurate the choice of  $\nu^{\text{chaos}}$  is, and the better the PCE identification is.

If another  $(N \times \nu^{\text{chaos},*})$  real matrix  $[\Psi^*]$  of independent realizations is considered, the matrices  $[B^*]$  and  $[\tilde{\Psi}^*] = [B^*]^{-1}[\Psi^*]$  can be computed according to Eqs. (3.33) and (3.34). As it has previously been seen,  $[\Psi^*]$ ,  $[\tilde{\Psi}]$  and  $[\tilde{\Psi}^*]$  keep the same structure. The accuracy of  $[z_{\mathcal{O}_\eta}^M]$  can thus be estimated by comparing  $\mathcal{C}([\nu^{\text{exp}}], [z_{\mathcal{O}_\eta}^M], [\Psi][B]^{-1})$  and  $\mathcal{C}([\nu^{\text{exp}}], [z_{\mathcal{O}_\eta}^M], [\Psi^*][B^*]^{-1})$ .

In particular,  $\nu^{\text{chaos},*}$  and  $\nu^{\text{chaos}}$  can be different. Finally, once the coefficient matrix  $[z_{\mathcal{O}_\eta}^M]$  has been computed, the higher  $\nu^{\text{chaos},*}$  is, the more accurate and general the validation is.

**4. Application.** The method proposed in the two former sections is applied to two examples that both deal with the identification of the truncated PCE coefficients of a random vector  $\boldsymbol{\eta}$  characterized by a multidimensional analytical distributions. The first one is built with  $N_\eta = 3$ , the second one with  $N_\eta = 50$ . According to the notations of the former sections, the set of the  $\nu^{\text{exp}}$  independent realizations used in the PCE identification are gathered in  $[\eta^{\text{exp}}]$ . Another set of  $\nu^{\text{ref}}$  independent realizations ( $\nu^{\text{ref}} \gg \nu^{\text{exp}}$ ) is used as a reference to validate the different modelings. The idea of this section is thus to show to what extent the whole method described in the former sections allows computing convergence analysis and relevant PCE identification of the random vector  $\boldsymbol{\eta}$  from a limited number of information,  $[\eta^{\text{exp}}]$ . Moreover, a distinction has to be made between the PDF modeling, achieved thanks to a PCE, and its estimation from PCE samples, computed thanks to non parametrical methods. In this context, let  $\nu^{\text{chaos}}$  be the number of independent realizations used to carry out the PCE identification, and  $\nu^{\text{chaos},*}$  the number of independent realizations of the identified PCE random vector, which will be used to draw graphical representations. For the two applications, the Gaussian measure is chosen for the PCE of  $\boldsymbol{\eta}$ , and  $\nu^{\text{exp}} = 1000$ .

It is reminded that, in this work, the term high dimension refers to the fact that the dimension  $N_\eta$  of unknown random vector  $\boldsymbol{\eta}$  is high.

**4.1. Application in low dimension.** The objective of this section is to apply the whole PCE method to a  $N_\eta = 3$ -dimension case. First, the statistical properties of the unknown random vector  $\boldsymbol{\eta}$  are presented. Secondly, a convergence analysis is carried out in order to calculate the optimal truncation parameters  $N_g$  and  $p$  of the PCE,  $\boldsymbol{\eta}^{\text{chaos}}(N)$ , of  $\boldsymbol{\eta}$ . Then, the PCE coefficients are identified from the  $\nu^{\text{exp}}$  independent realizations,  $[\eta^{\text{exp}}]$ , of  $\boldsymbol{\eta}$ . At last, the relevance of the PCE modeling is analysed.

**4.1.1. Generation of the random vector to identify.** Let  $[X]$  be a  $(3 \times 6)$  real-valued random matrix whose coefficients are uniformly and independently chosen between -1 and 1, such that  $\boldsymbol{\eta}$  is defined according to the notations of Section 3 as:

$$\boldsymbol{\eta} = [X] \mathcal{E}(\boldsymbol{\xi}^{\text{exp}}, 2), \quad (4.1)$$

where  $\boldsymbol{\xi}^{\text{exp}} = (\xi_1^{\text{exp}}, \xi_2^{\text{exp}})$  is a normalized Gaussian random vector which components are independent. The components of  $\boldsymbol{\eta}$  are however strongly dependent, and the PCE truncation parameters to be found back by the convergence analysis are  $p^{\text{exp}} = 2$  and  $N_g^{\text{exp}} = 2$ .

Let  $\{\boldsymbol{\xi}^{\text{exp}}(\theta_1), \dots, \boldsymbol{\xi}^{\text{exp}}(\theta_{\nu^{\text{exp}}})\}$  and  $\{\boldsymbol{\xi}^{\text{exp}}(\theta_1), \dots, \boldsymbol{\xi}^{\text{exp}}(\theta_{\nu^{\text{ref}}})\}$  be  $\nu^{\text{exp}}$  and  $\nu^{\text{ref}}$  independent realizations of the random vector  $\boldsymbol{\xi}^{\text{exp}}$ , such that the matrices of independent realizations  $[\eta^{\text{exp}}]$  and  $[\eta^{\text{ref}}]$  are given by:

$$[\eta^{\text{exp}}] = [X] [\mathcal{E}(\boldsymbol{\xi}^{\text{exp}}(\theta_1), 2) \quad \dots \quad \mathcal{E}(\boldsymbol{\xi}^{\text{exp}}(\theta_{\nu^{\text{exp}}}), 2)], \quad (4.2)$$

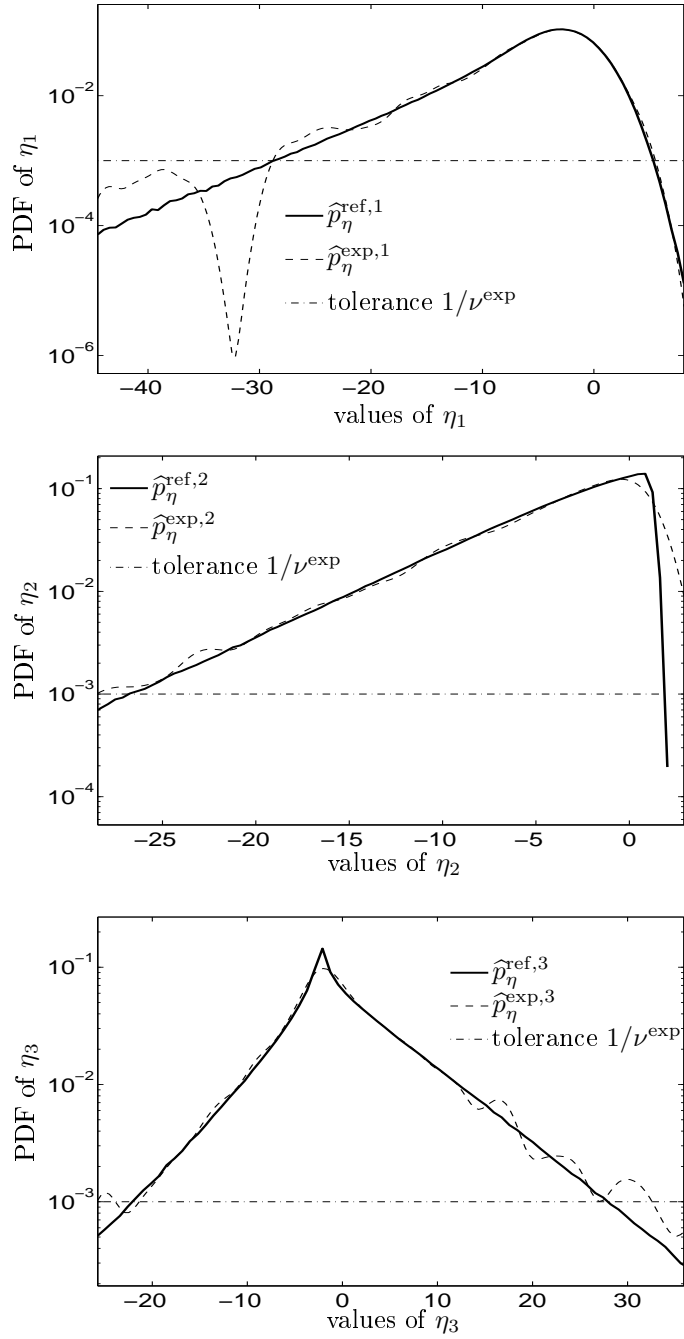
$$[\eta^{\text{ref}}] = [X] [\mathcal{E}(\boldsymbol{\xi}^{\text{exp}}(\theta_1), 2) \quad \dots \quad \mathcal{E}(\boldsymbol{\xi}^{\text{exp}}(\theta_{\nu^{\text{ref}}}), 2)]. \quad (4.3)$$

Let  $\{\widehat{p}_\eta^{\text{ref},k}, 1 \leq k \leq 3\}$  be the Kernel smoothing estimations of the marginal PDFs of each component of  $\boldsymbol{\eta}$ , which are computed thanks to the  $\nu^{\text{ref}}$  independent realizations of  $\boldsymbol{\eta}$  gathered in  $[\eta^{\text{ref}}]$ . In this example,  $\nu^{\text{ref}} = 2 \times 10^6 \gg \nu^{\text{exp}} = 1000$ . It is reminded that the PCE identification of  $\boldsymbol{\eta}$  is only achieved thanks to the matrix of independent realizations  $[\eta^{\text{exp}}]$ , which is considered as the only available information. The PDFs  $\{\widehat{p}_\eta^{\text{ref},k}, 1 \leq k \leq 3\}$  are moreover supposed to build the marginal PDFs of the reference  $\boldsymbol{\eta}$ .

**4.1.2. Identification of the PCE truncation parameters.** Using the notations of Section 2.4, the boundary intervals  $\text{BI}_1$ ,  $\text{BI}_2$  and  $\text{BI}_3$  for which the convergence analysis is achieved, are chosen such that:

$$\forall 1 \leq k \leq 3, \text{BI}_k = \left\{ x \in \mathbb{R} \mid \widehat{p}_\eta^{\text{ref},k}(x) \geq \frac{1}{\nu^{\text{exp}}} \right\}. \quad (4.4)$$

Figure 4.1 displays the reference marginal PDFs of  $\boldsymbol{\eta}$ , as well as the marginal PDFs  $\{\widehat{p}_\eta^{\text{exp},k}, 1 \leq k \leq 3\}$  estimated from the  $\nu^{\text{exp}}$  independent realizations only. The  $1/\nu^{\text{exp}}$  tolerance is also plotted so that the boundary intervals can therefore be deduced from these graphs.

Figure 4.1: Graphs of the marginal PDFs of  $\boldsymbol{\eta}$ .

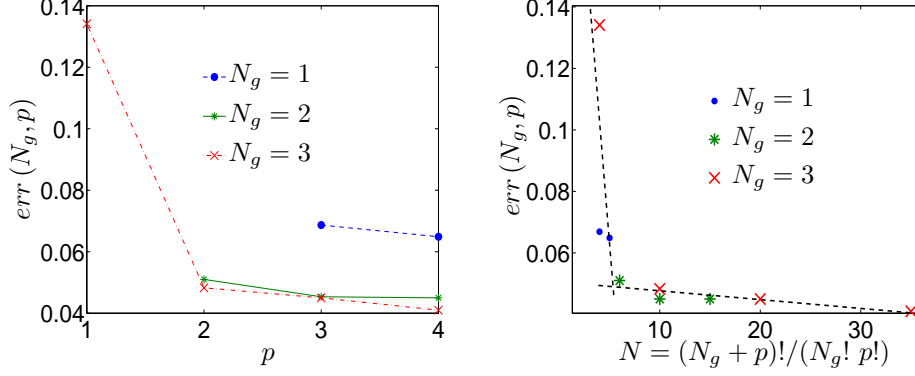
Figure 4.2: Convergence analysis of the PCE of  $\eta$ .

Figure 4.2 shows the values of  $err(N_g, p)$ , for nine pairs  $(N_g, p)$  in  $\mathcal{Q}(3)$ . On these graphs, the gradient break of  $N \mapsto err(N)$  is observed at  $N = 6$ , which allows us to find back the initial solution  $p^{\text{exp}} = 2$  and  $N_g^{\text{exp}} = 2$ . For this small dimension case, the optimal truncation parameters  $p$  and  $N_g$  given by the convergence analysis are equal to the parameters of the analytical reference PCE.

**4.1.3. PCE Identification.** The former convergence analysis leads us to the following PCE of  $\eta$ :

$$\eta \approx \eta^{\text{chaos}}(6) = \sum_{j=1}^6 \mathbf{y}^j \Psi_j(\xi_1, \xi_2) = [\mathbf{y}] \Psi(\xi_1, \xi_2), \quad (4.5)$$

where  $\xi_1$  and  $\xi_2$  are two independent normalized Gaussian random variables. We are now going to compare  $[\mathbf{y}^{\text{class}}]$  and  $[\mathbf{y}^{\text{new}}]$ , where  $[\mathbf{y}^{\text{class}}]$  stems from the classical problem defined by Eq. (2.24), whereas  $[\mathbf{y}^{\text{new}}]$  comes from the maximization of the new formulation defined by Eq. (3.40). In this application,  $\nu^{\text{chaos}} = 1000$ , and the two PCE identifications have been computed thanks to the same numerical cost  $M = 10^4$ , which means that  $M = 10^4$  independent random trials of  $[\mathbf{y}^{\text{class}}]$  and  $[\mathbf{y}^{\text{new}}]$  have been used to maximize the log-likelihood. The value of  $M$  has been chosen sufficiently high for the PCE error function  $err(N_g, p)$  to be independent of  $M$ . Hence, for a new matrix of independent realizations,  $[\Psi^*]$ , of size  $(6 \times \nu^{\text{chaos},*})$ , independent realizations  $[\eta^{\text{class}}(6)]$  and  $[\eta^{\text{new}}(6)]$  of  $\eta^{\text{chaos}}(6)$  are deduced, with respect to the two optimization options:

$$[\eta^{\text{class}}(6)] = [\mathbf{y}^{\text{class}}][\Psi^*], \quad (4.6)$$

$$[\eta^{\text{new}}(6)] = [\mathbf{y}^{\text{new}}][\Psi^*]. \quad (4.7)$$

Let

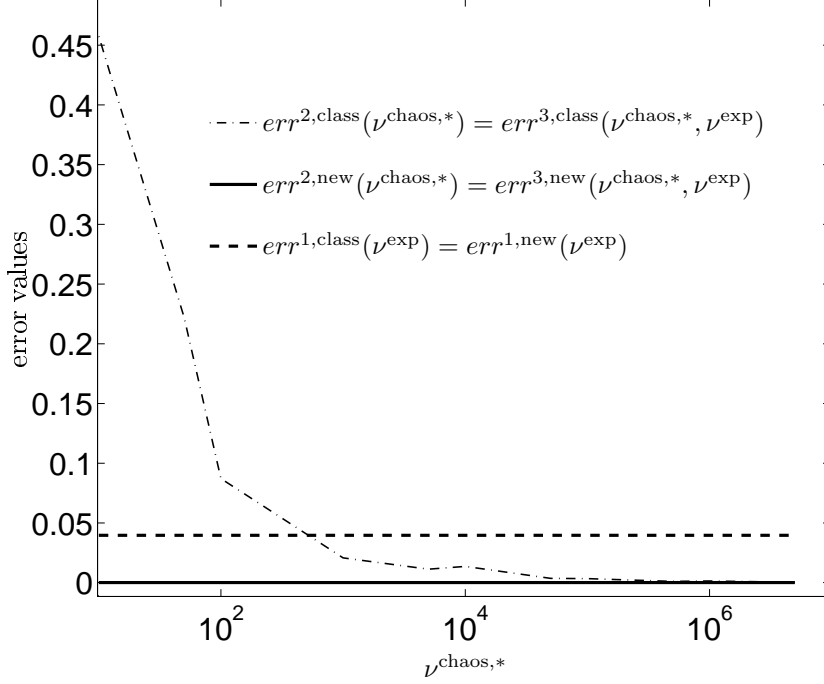


Figure 4.3: Convergence of the autocorrelation error functions with respect to  $\nu^{\text{chaos},*}$ .

$$[R_{\eta}^{\text{exp}}] = \frac{1}{\nu^{\text{exp}}} [\eta^{\text{exp}}] [\eta^{\text{exp}}]^T, \quad (4.8)$$

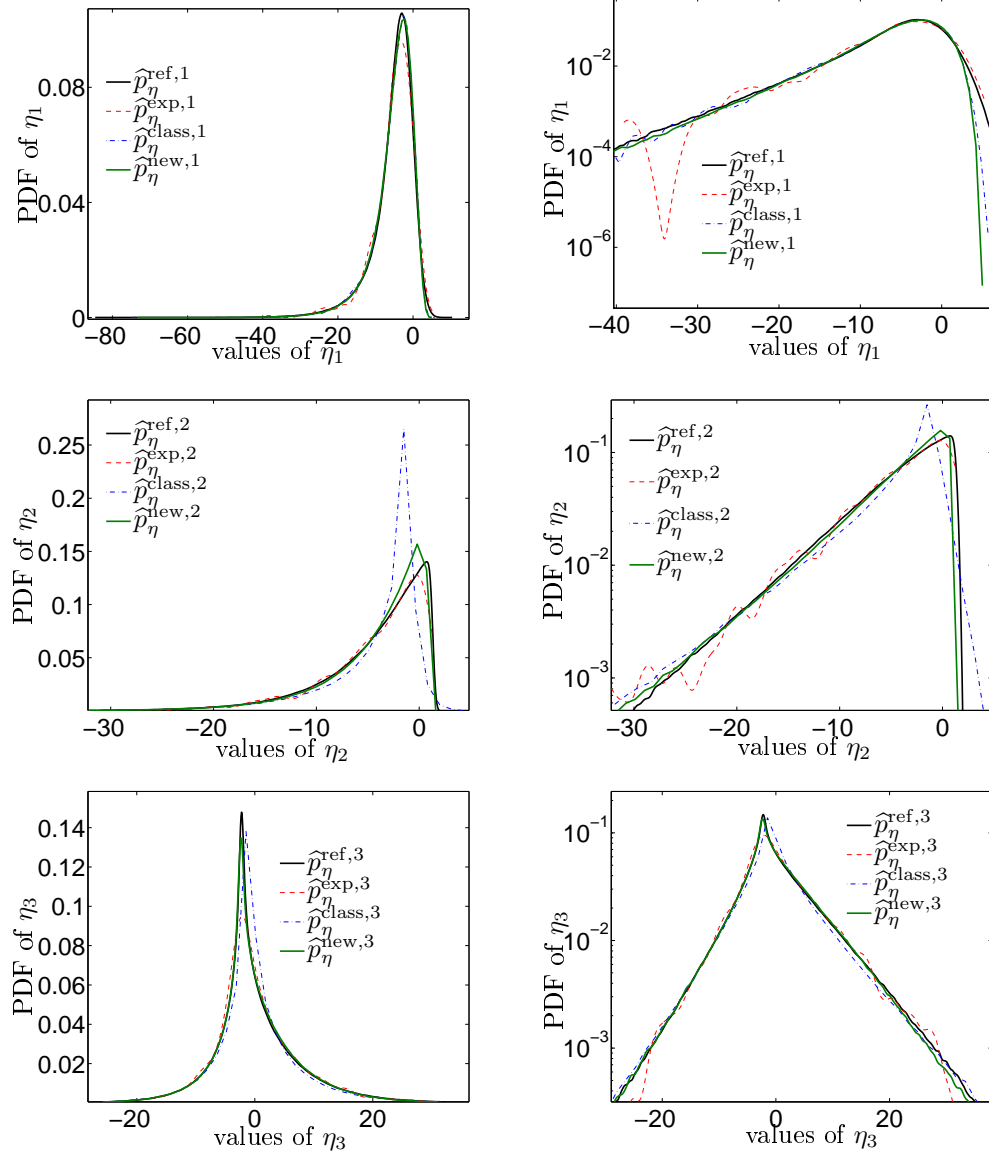
$$[R_{\eta}^{\text{ref}}] = \frac{1}{\nu^{\text{ref}}} [\eta^{\text{ref}}] [\eta^{\text{ref}}]^T, \quad (4.9)$$

$$[R_{\eta}^{\text{class}}] = \frac{1}{\nu^{\text{chaos},*}} [\eta^{\text{class}}(6)] [\eta^{\text{class}}(6)]^T, \quad (4.10)$$

$$[R_{\eta}^{\text{new}}] = \frac{1}{\nu^{\text{chaos},*}} [\eta^{\text{new}}(6)] [\eta^{\text{new}}(6)]^T \quad (4.11)$$

be four estimations of the autocorrelation matrix  $[R_{\eta}]$  of  $\boldsymbol{\eta}$ . It is supposed that  $[R_{\eta}^{\text{ref}}]$  is the best approximation of  $[R_{\eta}]$  and will be considered as the reference. According to the Eqs. (3.4), (3.5) and (3.6), the autocorrelation errors  $err^{1,\text{class}}$ ,  $err^{2,\text{class}}$ ,  $err^{3,\text{class}}$  and  $err^{1,\text{new}}$ ,  $err^{2,\text{new}}$ ,  $err^{3,\text{new}}$  are then computed in each case. In figure 4.3, it can thus be verified that:

$$\forall \nu^{\text{chaos},*} \geq 6, \quad err^{2,\text{new}}(\nu^{\text{chaos},*}) = err^{3,\text{new}}(\nu^{\text{chaos},*}, \nu^{\text{exp}}) = 0, \quad (4.12)$$

Figure 4.4: Comparison of the marginal PDFs of  $\eta$  and  $\eta^{\text{chaos}(6)}$ .

$$\lim_{\nu^{\text{chaos},*} \rightarrow +\infty} \text{err}^{2,\text{class}}(\nu^{\text{chaos},*}) = \text{err}^{3,\text{class}}(\nu^{\text{chaos},*}, \nu^{\text{exp}}) = 0. \quad (4.13)$$

In particular, for the value  $\nu^{\text{chaos},*} = \nu^{\text{chaos}} = 1000$ , it can be noticed that the values of  $\text{err}^{2,\text{class}}(\nu^{\text{chaos},*})$  and  $\text{err}^{3,\text{class}}(\nu^{\text{chaos},*}, \nu^{\text{exp}})$  are significant when compared to  $\text{err}^{1,\text{class}}(\nu^{\text{exp}})$ , which introduces an additive bias in the identification.



Figure 4.4 shows a comparison between the marginal PDFs  $\hat{p}_\eta^{\text{ref},k}$ ,  $\hat{p}_\eta^{\text{exp},k}$ ,  $\hat{p}_\eta^{\text{class},k}$  and  $\hat{p}_\eta^{\text{new},k}$ , for  $1 \leq k \leq 3$ . These PDFs are estimated using Kernel smoothing on the independent realizations gathered in the matrices  $[\eta^{\text{ref}}]$ ,  $[\eta^{\text{exp}}]$ ,  $[\eta^{\text{class}}(6)]$  and  $[\eta^{\text{new}}(6)]$ , respectively, with  $\nu^{\text{chaos},*} = 10^6 \gg \nu^{\text{chaos}} = 1000$ . First, from only  $\nu^{\text{exp}} = 1000$  independent realizations of  $\boldsymbol{\eta}$ , it can be seen that the marginal PDFs are well described by the PCE random vectors  $\boldsymbol{\eta}^{\text{new}}(6)$  and  $\boldsymbol{\eta}^{\text{class}}(6)$ . In particular, the PDFs tails are very well characterized. The PCE method is therefore an extremely efficient tool to build arbitrary multidimensional PDFs. Secondly, it can be noticed that, for a same computational cost  $M$ , the new PCE identification formulation leads us to better results than the classical one. Finally, to still improve these PCE, more trials in  $\mathcal{O}_\eta$  would be necessary to better characterize  $[y^{\text{class}}]$  and  $[y^{\text{new}}]$ . In order to obtain a PCE that corresponds still more precisely to the reference random vector  $\boldsymbol{\eta}$ , an increase of  $\nu^{\text{exp}}$ , that is to say, more information about  $\boldsymbol{\eta}$ , would have been required.

#### 4.1.4. Relevance of the PCE compared to Kernel Mixture and PASM.

From adequacy tests, likelihood estimations and graphical representations, the idea of this section is to show the assets of the new PCE formulation when dealing with the identification of multidimensional distributions from a limited knowledge on the random vector of interest  $\boldsymbol{\eta}$  compared to Kernel Mixture (KM) and Prior Algebraic Stochastic Modeling (PASM). In this prospect, two PDFs  $\hat{p}_\eta(\mathbf{x})$  and  $\hat{p}_\eta^{\text{PASM}}(\mathbf{x}, \mathbf{w})$  are built using a KM approach and a PASM method. The input data of these modelings are still the matrix of independent realizations  $[\eta^{\text{exp}}] = [\boldsymbol{\eta}^{(1)} \ \dots \ \boldsymbol{\eta}^{(\nu^{\text{exp}})}]$ . Once the KM, the PASM and the two PCE projection matrices,  $[y^{\text{class}}]$  and  $[y^{\text{new}}]$ , are constructed,  $Q$  independent realizations are computed from the four distributions, from which comparisons to the reference solution are achieved. For this application,  $Q = 10^6$ .

#### Construction of independent realisations.

##### • Kernel Mixture.

Considering an independent Gaussian multidimensional Kernel, a non parametrical PDF  $\hat{p}_\eta(\mathbf{x})$  is postulated as a sum of  $\nu^{\text{exp}}$  Gaussian PDFs  $\{p_i, 1 \leq i \leq \nu^{\text{exp}}\}$  to model  $p_\eta(\mathbf{x})$ :

$$\hat{p}_\eta(\mathbf{x}) = \sum_{i=1}^{\nu^{\text{exp}}} \frac{1}{\nu^{\text{exp}}} p_i(\mathbf{x}), \quad (4.14)$$

$$p_i(\mathbf{x}) = \prod_{k=1}^{N_\eta} \frac{1}{\sqrt{2\pi}h_k} \exp\left(-\frac{1}{2}\left(\frac{x_k - \eta_k^i}{h_k}\right)^2\right), \quad (4.15)$$

$$\mathbf{h} = \hat{\sigma} \left( \frac{4}{(2 + N_\eta) \nu^{\text{exp}}} \right)^{1/(N_\eta+4)}, \quad (4.16)$$

where  $\mathbf{x} \mapsto p_i(\mathbf{x})$  is the  $N_\eta$ -dimension multivariate Gaussian PDF, with mean value

$$\boldsymbol{\eta}^{(i)} \text{ and covariance matrix } \begin{pmatrix} h_1 & 0 & \dots & 0 \\ 0 & h_2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & h_{N_\eta} \end{pmatrix}, \mathbf{h} \text{ is the multidimensionnal op-}$$

timal Silverman bandwidth, and  $\widehat{\sigma}_k$  is the empirical estimation of the standard deviation of each component  $\eta_k$  of  $\boldsymbol{\eta}$ . Let  $\boldsymbol{\eta}^{\text{ker}}$  be the Kernel Mixture characterized by the PDF  $\boldsymbol{x} \mapsto \widehat{p}_{\boldsymbol{\eta}}(\boldsymbol{x})$ . The  $Q$  independent realizations  $\{\boldsymbol{\eta}^{\text{ker},1}, \dots, \boldsymbol{\eta}^{\text{ker},Q}\}$  of  $\boldsymbol{\eta}^{\text{ker}}$  are then computed and gathered in the matrix  $[\boldsymbol{\eta}^{\text{ker}}]$ .

• **Prior Algebraic Stochastic Modeling.**

From the  $\nu^{\text{exp}}$  independent realizations of  $\boldsymbol{\eta}$ , the  $N_{\eta}$  marginal cumulative distributions  $F_{\eta_k}$  of  $\eta_k$ , with  $1 \leq k \leq N_{\eta}$ , are estimated using a non parametric statistical method. In addition, a Gaussian copula  $C_{\text{rank}}^{\text{gauss}}$  (see [6] for more details about the copula) based on the rank correlation is chosen (this type of copula has been chosen as it is the most commonly used in the PASM approaches):

$$C_{\text{rank}}^{\text{gauss}}(x_1, \dots, x_{N_{\eta}}) = \phi_{\text{rank}}^{N_{\eta}}(\phi^{-1}(x_1), \dots, \phi^{-1}(x_{N_{\eta}})), \quad (4.17)$$

$$\phi_{\text{rank}}^{N_{\eta}}(\mathbf{u}) = \int_{-\infty}^{u_1} \dots \int_{-\infty}^{u_{N_{\eta}}} \frac{1}{(2\pi)^{N_{\eta}/2} \sqrt{\det([R^{\text{rank}}])}} \exp\left(-\frac{1}{2}\mathbf{u}^T [R^{\text{rank}}] \mathbf{u}\right) du_1 \dots du_{N_{\eta}}, \quad (4.18)$$

$$\phi(v) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^v \exp\left(-\frac{v^2}{2}\right) dv, \quad (4.19)$$

$$[R^{\text{rank}}]_{ij} = 2 \sin\left(\frac{\pi}{6} \rho_{ij}^S\right), \quad (4.20)$$

where  $\rho_{ij}^S$  is the Spearman correlation coefficient between  $\eta_i$  and  $\eta_j$ . Let  $\boldsymbol{\eta}^{\text{cop}}$  be the random vector characterized by the copula  $C_{\text{rank}}^{\text{gauss}}$  and the marginal cumulative distributions  $\{F_{\eta_k}, 1 \leq k \leq N_{\eta}\}$ .  $Q$  independent realizations of  $\boldsymbol{\eta}^{\text{cop}}$  are thus gathered in the matrix  $[\boldsymbol{\eta}^{\text{cop}}]$ .

• **Polynomial chaos expansion.**

Finally, using the matrices  $[y^{\text{class}}]$  and  $[y^{\text{new}}]$  of Section 4.1.3, and a new  $(6 \times Q)$  real matrix  $[\Psi^Q]$  of realizations,  $Q$  independent realizations of  $\boldsymbol{\eta}^{\text{class}}(6)$  and  $\boldsymbol{\eta}^{\text{new}}(6)$  are gathered in the matrix  $[\boldsymbol{\eta}^{\text{class}}] = [y^{\text{class}}][\Psi^Q]$  and  $[\boldsymbol{\eta}^{\text{new}}] = [y^{\text{new}}][\Psi^Q]$ .

**Relevance of the PCE modeling when identifying multidimensional PDFs from a limited amount of independent realizations.** Using the results of Parametrical Statistics, this section assesses the relevance of the four methods to construct multidimensional PDFs. Three kinds of analysis are achieved: adequacy tests, 2D graphical representations, and multidimensional likelihood computations.

• **Adequacy tests.**

From the matrices of independent realizations  $[\boldsymbol{\eta}^{\text{ker}}]$ ,  $[\boldsymbol{\eta}^{\text{cop}}]$ ,  $[\boldsymbol{\eta}^{\text{class}}]$  and  $[\boldsymbol{\eta}^{\text{new}}]$ , the estimations  $\{\widehat{F}_k^{\text{ker}}, 1 \leq k \leq N_{\eta}\}$ ,  $\{\widehat{F}_k^{\text{cop}}, 1 \leq k \leq N_{\eta}\}$ ,  $\{\widehat{F}_k^{\text{class}}, 1 \leq k \leq N_{\eta}\}$  and  $\{\widehat{F}_k^{\text{new}}, 1 \leq k \leq N_{\eta}\}$  of the cumulative distribution functions (CDF) of each components of  $\boldsymbol{\eta}^{\text{ker}}$ ,  $\boldsymbol{\eta}^{\text{cop}}$ ,  $\boldsymbol{\eta}^{\text{class}}(6)$  and  $\boldsymbol{\eta}^{\text{new}}(6)$  are respectively assessed. Let  $\tilde{\boldsymbol{\eta}}^{(1)}, \dots, \tilde{\boldsymbol{\eta}}^{(N_{\eta})}$  be the  $1 \times \nu^{\text{exp}}$ -dimension linear forms corresponding to the lines of  $[\boldsymbol{\eta}^{\text{exp}}]$ . For

$1 \leq k \leq N_\eta$ ,  $\tilde{\eta}^{(k)}$  gathers therefore the  $\nu^{\text{exp}}$  independent realizations of the component  $\eta_k$  of  $\eta$ , which have been used to compute the statistical modelings. For  $1 \leq k \leq N_\eta$ , the Kolmogorov-Smirnov adequacy tests are then performed. For each component  $\eta_k$  of  $\eta$ , the null distribution of the Kolmogorov-Smirnov statistics is computed under the null hypothesis that the  $\nu^{\text{exp}}$  independent realizations of  $\tilde{\eta}^{(k)}$  are drawn from the distribution of the chosen stochastic model. Table 4.1 gives the  $\beta$ -value for each stochastic model, which is defined as the probability of obtaining a test statistic at least as extreme as the one that was actually observed, assuming that the null hypothesis is true. Without surprise, this table allows us to verify that the modeling based on the Gaussian copula and the empirical PDFs of each components of  $\eta$  gives the best results. Moreover, with an error level of 5%, only the tests for the copula model and the PCE identification based on the new formulation are positive. The classical PCE and the Kernel mixture modelings are indeed less relevant to characterize the marginal PDFs of  $\eta$ .

CDF	$\hat{F}_1^{\text{class}}$	$\hat{F}_1^{\text{new}}$	$\hat{F}_1^{\text{ker}}$	$\hat{F}_1^{\text{cop}}$
$\beta$ -value	0.3779	0.6331	0.2142	0.9996
CDF	$\hat{F}_2^{\text{class}}$	$\hat{F}_2^{\text{new}}$	$\hat{F}_2^{\text{ker}}$	$\hat{F}_2^{\text{cop}}$
$\beta$ -value	0.0000	0.0967	0.0000	0.4573
CDF	$\hat{F}_3^{\text{class}}$	$\hat{F}_3^{\text{new}}$	$\hat{F}_3^{\text{ker}}$	$\hat{F}_3^{\text{cop}}$
$\beta$ -value	0.0000	0.8692	0.0411	0.9849

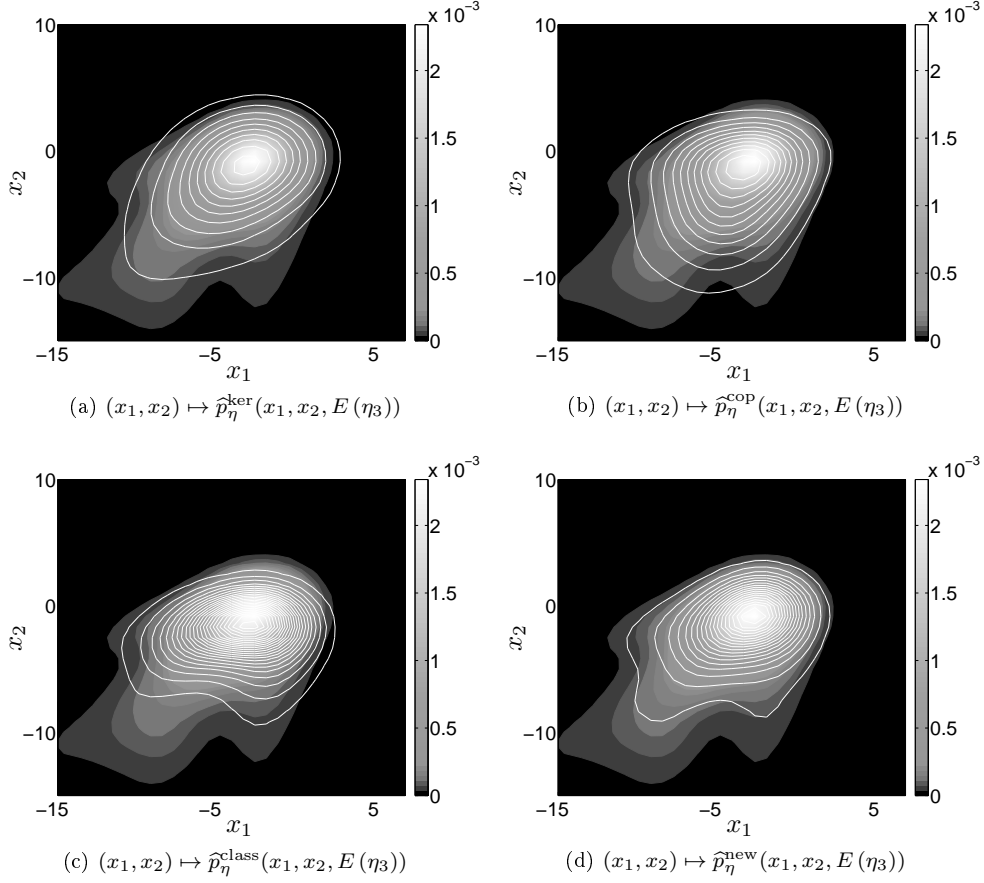
Table 4.1: Computation of the  $\beta$ -values corresponding to the different stochastic models.

- **Two-dimensions graphical analysis.**

From  $[\eta^{\text{ref}}]$ ,  $[\eta^{\text{ker}}]$ ,  $[\eta^{\text{cop}}]$ ,  $[\eta^{\text{class}}]$  and  $[\eta^{\text{new}}]$ , the estimations  $\mathbf{x} \mapsto \hat{p}_\eta^{\text{ref}}(\mathbf{x})$ ,  $\mathbf{x} \mapsto \hat{p}_\eta^{\text{ker}}(\mathbf{x})$ ,  $\mathbf{x} \mapsto \hat{p}_\eta^{\text{cop}}(\mathbf{x})$ ,  $\mathbf{x} \mapsto \hat{p}_\eta^{\text{class}}(\mathbf{x})$  and  $\mathbf{x} \mapsto \hat{p}_\eta^{\text{new}}(\mathbf{x})$  of the multidimensional PDF of  $\eta$ ,  $\eta^{\text{ker}}$ ,  $\eta^{\text{cop}}$ ,  $\eta^{\text{class}}$ (6),  $\eta^{\text{new}}$ (6) are respectively computed using the non parametric statistical estimation defined by Eq. (2.18). Projections of these functions are presented in Figures 4.5, 4.6 and 4.7. In each figure, the surface plot characterizes the reference PDF (based on the  $\nu^{\text{ref}} = 2 \times 10^6$  independent realizations), and the contour plot refers to isovalues of the projected PDF of interest. It can therefore be seen that the new formulation of the PCE gives very good results in identifying multidimensional PDFs. In addition, in this example, the Kernel mixture model is more adapted than the copula based model to characterize the multidimensional PDFs.

- **Likelihood estimations.**

From Eq. (2.10), the multidimensional log-likelihood functions  $\mathcal{L}_{\eta^{\text{ker}}}([\eta^{\text{exp}}])$ ,  $\mathcal{L}_{\eta^{\text{cop}}}([\eta^{\text{exp}}])$ ,  $\mathcal{L}_{\eta^{\text{class}}}([\eta^{\text{exp}}])$  and  $\mathcal{L}_{\eta^{\text{new}}}([\eta^{\text{exp}}])$  are estimated from the realizations matrices  $[\eta^{\text{exp}}]$ ,  $[\eta^{\text{ker}}]$ ,  $[\eta^{\text{cop}}]$ ,  $[\eta^{\text{class}}]$  and  $[\eta^{\text{new}}]$ , in order to evaluate the multidimensional relevance of the different stochastic models. In the same manner,  $[\eta^{\text{ref}}]_{1000}$  is defined as the 1000 first columns of  $[\eta^{\text{ref}}]$ , and the log-likelihood functions  $\mathcal{L}_{\eta^{\text{ker}}}([\eta^{\text{ref}}]_{1000})$ ,  $\mathcal{L}_{\eta^{\text{cop}}}([\eta^{\text{ref}}]_{1000})$ ,  $\mathcal{L}_{\eta^{\text{class}}}([\eta^{\text{ref}}]_{1000})$  and  $\mathcal{L}_{\eta^{\text{new}}}([\eta^{\text{ref}}]_{1000})$  are computed. These values are gathered in Table 4.2. It can thus be verified that the new formulation of the PCE identification gives the best results when considering the maximization of the log-likelihood.

Figure 4.5: Comparison of 2D contours plots in the plane  $[x_3 = E(\eta_3)]$ .

$\mathcal{L}_{\eta^{\text{ker}}}([\eta^{\text{exp}}])$	$\mathcal{L}_{\eta^{\text{cop}}}([\eta^{\text{exp}}])$	$\mathcal{L}_{\eta^{\text{class}}}([\eta^{\text{exp}}])$	$\mathcal{L}_{\eta^{\text{new}}}([\eta^{\text{exp}}])$
$-8.0712 \cdot 10^3$	$-8.7530 \cdot 10^3$	$-8.1844 \cdot 10^3$	$-7.8624 \cdot 10^3$
$\mathcal{L}_{\eta^{\text{ker}}}([\eta^{\text{ref}}]_{1000})$	$\mathcal{L}_{\eta^{\text{cop}}}([\eta^{\text{ref}}]_{1000})$	$\mathcal{L}_{\eta^{\text{class}}}([\eta^{\text{ref}}]_{1000})$	$\mathcal{L}_{\eta^{\text{new}}}([\eta^{\text{ref}}]_{1000})$
$-8.1933 \cdot 10^3$	$-8.5535 \cdot 10^3$	$-8.1797 \cdot 10^3$	$-7.8457 \cdot 10^3$

Table 4.2: Computation of the multidimensional log-likelihood corresponding to the different stochastic models.

As a conclusion for this example, in low dimension, it can be seen that the new formulation of the PCE identification is very relevant when trying to identify multidimensional distributions from a limited number of measurements. Indeed, it allows us to build multidimensional distributions that are still relevant for experimental data that have not been used in the identification process.

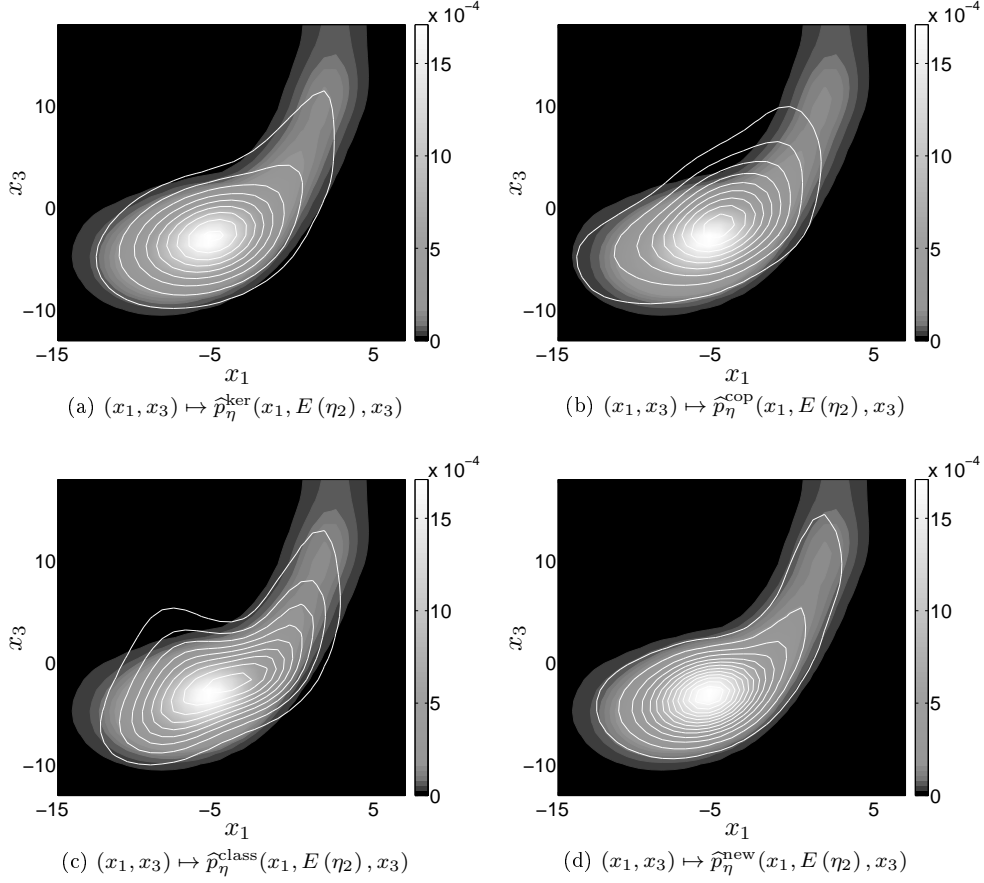
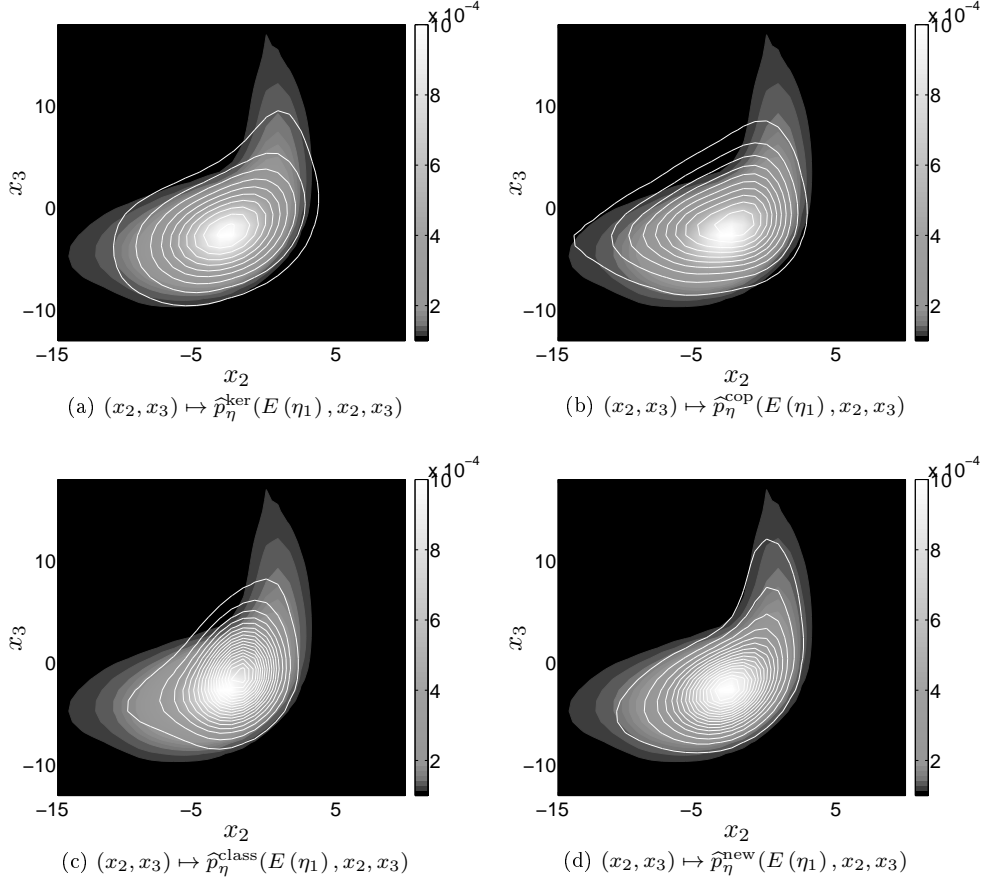


Figure 4.6: Comparison of 2D contours plots in the plane  $[x_2 = E(\eta_2)]$ .

**4.2. Application in high dimension.** The idea of this second application is to underline the ability of the new PCE formulation to carry out convergence analysis in high dimension. Indeed, as it has been shown in Section 3, for a given value of  $\nu^{\text{chaos}}$ , when the size  $N$  of the polynomial basis increases, and more specially when the maximum degree  $p$  of the polynomial basis becomes high, the difference  $\frac{1}{\nu^{\text{chaos}}}[\Psi][\Psi]^T - [I_N]$  introduces a significant numerical bias which perturbs the classical PCE identification. In opposite, the new PCE formulation, which avoids computational autocorrelation errors, allows the numerical algorithms to be much more stable and to give more relevant results.

**4.2.1. Generation of a high dimension random vector.** Using the same notations than in Sections 2.1 and 4.1.1, let  $[X^{\text{HD}}]$  be a  $(N_\eta \times N)$  real matrix whose entries are randomly generated, such that random vector  $\boldsymbol{\eta}$  is given by:

$$\boldsymbol{\eta} = [X^{\text{HD}}]\boldsymbol{\Psi}(\boldsymbol{\xi}^{\text{exp}}, p^{\text{exp}}), \quad (4.21)$$

Figure 4.7: Comparison of 2D contours plots in the plane  $[x_1 = E(\eta_1)]$ .

$$\boldsymbol{\xi}^{\text{exp}} = \left( \xi_1^{\text{exp}}, \xi_2^{\text{exp}}, \dots, \xi_{N_g}^{\text{exp}} \right), \quad (4.22)$$

where  $\{\xi_k^{\text{exp}}, 1 \leq k \leq N_\eta\}$  is a set of  $N_\eta$  independent normalized Gaussian random variables. As in Section 4.1, we define a  $(N_\eta \times \nu^{\text{exp}})$  real matrix  $[\eta^{\text{exp}}]$ , which gathers  $\nu^{\text{exp}}$  independent realizations of  $\boldsymbol{\eta}$ :

$$[\eta^{\text{exp}}] = [X^{\text{HD}}][\Psi^{\text{exp}}], \quad (4.23)$$

$$[\Psi^{\text{exp}}] = [\Psi(\boldsymbol{\xi}^{\text{exp}}(\theta_1), p^{\text{exp}}) \quad \dots \quad \Psi(\boldsymbol{\xi}^{\text{exp}}(\theta_{\nu^{\text{exp}}}), p^{\text{exp}})]. \quad (4.24)$$

The components of the random vector  $\boldsymbol{\eta}$  are again strongly dependent. As a numerical illustration, it is supposed that  $\nu^{\text{exp}} = 1000$ ,  $p^{\text{exp}} = 9$ ,  $N_g^{\text{exp}} = 3$ ,  $N = (9+3)!/(9!3!) = 220$ ,  $N_\eta = 50$ . A high value of  $p^{\text{exp}}$  is deliberately chosen, in order to emphasize the

difficulties of the classical PCE formulation to carry out convergence analysis in high dimension. Nevertheless, this high value of  $p$  implies an ill-conditionning of  $[\Psi^{\text{exp}}]$ , such that  $\boldsymbol{\eta}$  can have very high values.

**4.2.2. Identification of the PCE truncation parameters.** According to Eq. (2.3), the truncated PCE,  $\boldsymbol{\eta}^{\text{chaos}}(N)$ , of  $\boldsymbol{\eta}$  is given by:

$$\boldsymbol{\eta}^{\text{chaos}}(N) = [y]\boldsymbol{\Psi}(\boldsymbol{\xi}, p). \quad (4.25)$$

Eq. (2.29) implies that the number  $N_y$  of elements of  $[y]$  has to be higher than  $N_\eta(N_\eta + 1)$ . When  $N_\eta$  is large, this leads us to the identification of thousands of coefficients. However, as it has been said in Section 2.4, the higher  $N_y$  is, the less precise is the PCE identification, for a given computational cost  $M$ . In addition, Section 3 has emphasized the ill-conditionning of matrix  $[\Psi]$  for high values of  $p$ . This motivates the definition of a new set  $\tilde{\mathcal{Q}}(p^{\text{max}}, N^{\text{max}})$ , such that the optimal values  $p^{\text{opt}}$  and  $N_g^{\text{opt}}$  are given by:

$$\tilde{\mathcal{Q}}(p^{\text{max}}, N^{\text{max}}) = \{(p, N_g), N_g \leq N_\eta, p \leq p^{\text{max}}, (N_g + p)! / (N_g! p!) \leq N^{\text{max}}\}, \quad (4.26)$$

$$(p^{\text{opt}}, N_g^{\text{opt}}) = \arg \min_{(p, N_g) \in \tilde{\mathcal{Q}}(p^{\text{max}}, N^{\text{max}})} \text{err}(N_g, p), \quad (4.27)$$

where  $\text{err}(N_g, p)$  is computed from  $M$  independent matrices in  $\mathcal{O}_\eta$ . For a fixed value  $\nu^{\text{chaos}} = 1000$ , the detrimental influence of the autocorrelation errors  $\text{err}^2$  and  $\text{err}^3$  of Eqs. (3.5) and (3.6) can then be noticed in Figure 4.8, when high values of  $N$  (and more specially high values of  $p$ ) are considered. The error functions  $\text{err}^{\text{class}}(N_g, p)$  and  $\text{err}^{\text{new}}(N_g, p)$  correspond, respectively, to the classical formulation and the new formulation of the PCE identification. It can be seen that for  $p \geq 8$ , the ratio  $\text{err}^{\text{class}}(N_g, p) / \text{err}^{\text{new}}(N_g, p)$  becomes greater than five, whereas the two methodologies are globally similar for low values of  $p$ . Hence, the accuracy of the classical method seems to be limited to low values of  $p$  and is therefore less relevant for convergence analysis which handle high polynomial orders. At last, the five lowest values of the numerical assessments of  $\text{err}^{\text{new}}(N_g, p)$  are gathered in Table 4.3. It can be seen that the new formulation allows finding back the couple  $(p^{\text{exp}}, N_g^{\text{exp}})$  as the minimum of the error function. Nevertheless, keeping in mind that the lowest  $N$  is, the easiest the identification is, this result also shows that using the couple  $(p, N_g) = (11, 2)$  could be interesting.

couples $(p, N_g)$	(11,2)	(9,3)	(7,4)	(6,5)	(2,27)
values of N	78	220	330	462	406
$\text{err}^{\text{new}}(N_g, p)$	0.06104	0.06005	0.06228	0.06301	0.06521

Table 4.3: Lowest values of  $\text{err}^{\text{new}}(N_g, p)$  with respect to  $(p, N_g)$ .

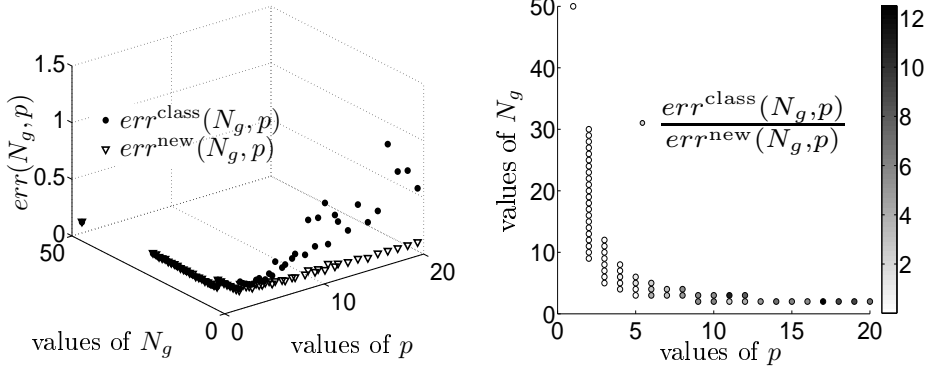


Figure 4.8: Comparison of the results for the convergence analysis of the two PCE identification formulations.

**4.3. PCE Identification.** From the  $\nu^{\text{exp}}$  independent realizations of  $\boldsymbol{\eta}$ , a PCE identification using the new formulation can be computed for the truncation parameters  $p = 9$  and  $N_g = 3$ , which correspond to  $N = 220$ . The results of the numerical identification with a computational cost of  $M = 1000$  are given in Figure 4.9. The value of  $M$  has been chosen for the PCE error function  $err(N_g, p)$  to be independent of  $M$ . In this figure, the marginal PDFs  $\hat{p}_{\eta_{41}}^{\text{chaos}}$  and  $\hat{p}_{\eta_{39}}^{\text{chaos}}$  of  $\eta_{41}^{\text{new}}(220)$  and  $\eta_{39}^{\text{new}}(220)$  are compared to the experimental estimations  $\hat{p}_{\eta_{41}}^{\text{exp}}$  and  $\hat{p}_{\eta_{39}}^{\text{exp}}$  of the components  $\eta_{41}$  and  $\eta_{39}$ , respectively. The values  $\eta_{41}^{\text{new}}(220)$  and  $\eta_{39}^{\text{new}}(220)$  correspond to the minimum and to the maximum values of the unidimensional error function  $err_k(3, 11)$ , for  $1 \leq k \leq 50$ , which is defined by Eq. (2.25). In order to evaluate the distance between these estimations and the true marginal PDFs of  $\boldsymbol{\eta}$ , the marginal PDFs estimated by the non parametric statistical Kernel method, with  $\nu^{\text{ref}} = 2 \times 10^5$  independent realizations of  $\eta_{41}$  and  $\eta_{39}$ , are added to the figures. These PDFs are considered as the reference. These figures therefore emphasize that the new PCE identification method allows building a stochastic model of the distribution of  $\boldsymbol{\eta}$  that suits the experimental marginal PDFs.

**5. Conclusion.** In the last decade, the increasing computational power has encouraged the development of computational models with increasing degrees of freedom. Hence, developing computational methods which can be applied to very high dimension cases is currently of great interest.

In this concern, this paper emphasized the efficiency of the PCE when building multidimensional distributions. After having quantified the detrimental influence of a numerical bias in the usual PCE identification methods in high dimension, this paper proposed a new formulation to allow performing relevant convergence analysis and PCE identification with respect to an arbitrary measure for the high dimension case. Finally, the method proposed allows making the PCE range reachable for many engineering applications with many degrees of freedom.

**Acknowledgments.** This work was supported by SNCF (Innovation and Research Department) and by the French Research Agency (Grant No: ANR-2010-BLAN-0904).



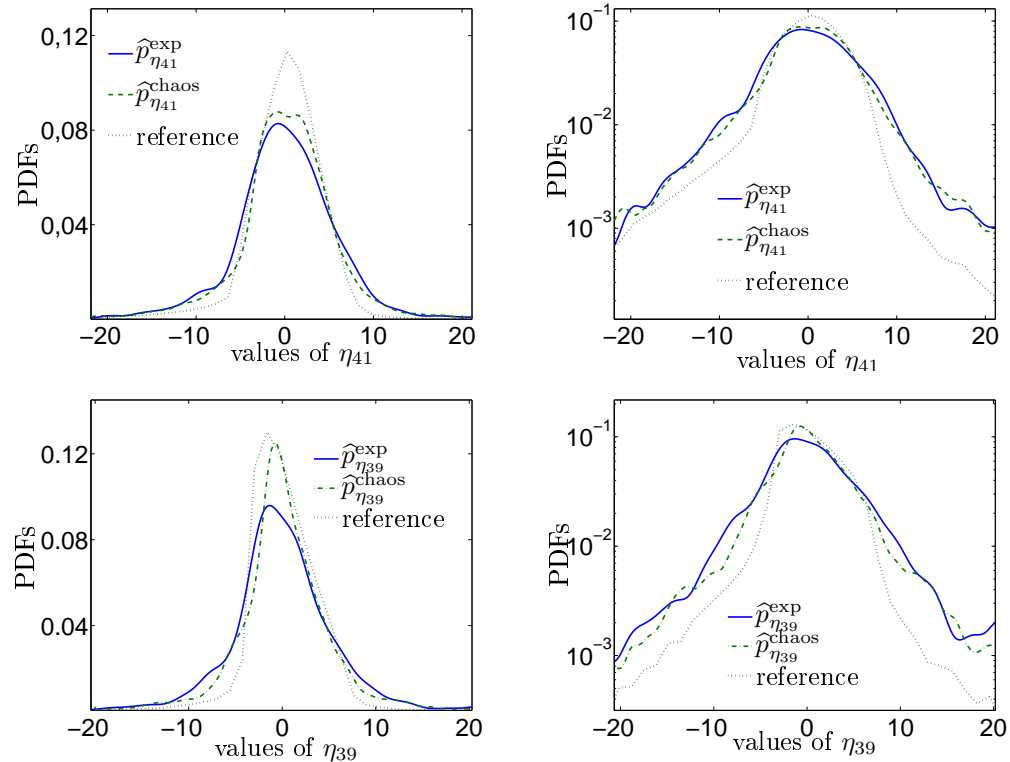


Figure 4.9: Graphs of the estimated marginal PDFs for two particular components of  $\eta$ .

#### REFERENCES

- [1] M. ARNST, R. GHANEM, AND C. SOIZE, *Identification of bayesian posteriors for coefficients of chaos expansions*, Journal of Computational Physics, 229 (9) (2010), pp. 3134–3154.
- [2] S. DAS, R. GHANEM, AND S. FINETTE, *Polynomial chaos representation of spatio-temporal random field from experimental measurements*, J. Comput. Phys., 228 (2009), pp. 8726–8751.
- [3] B. DEBUSSCHERE, H. NAJM, P. PÉBAY, O. M. KNIO, R. G. GHANEM, AND O. P. L. MAÎTRE, *Numerical challenges in the use of polynomial chaos representations for stochastic processes*, SIAM J. Sci. Comput., 26 (2004), pp. 698–719.
- [4] C. DESCIELERS, R. GHANEM, AND C. SOIZE, *Maximum likelihood estimation of stochastic chaos representations from experimental data*, Internat. J. Numer. Methods Engrg., 66 (2006), pp. 978–1001.
- [5] C. DESCIELERS, C. SOIZE, AND R. GHANEM, *Identification of chaos representations of elastic properties of random media using experimental vibration tests*, Comput. Mech., 39 (2007), pp. 831–838.
- [6] A. DUTFOY, *Reference guide*, tutorial, Open TURNS version 0.11.3, 2008.
- [7] R. GHANEM AND R. DOOSTAN, *Characterization of stochastic system parameters from experimental data: A bayesian inference approach*, J. Comput. Phys., 217 (2006), pp. 63–81.
- [8] R. GHANEM, S. MASRI, M. PELLISSETTI, AND R. WOLFE, *Identification and prediction of stochastic dynamical systems in a polynomial chaos basis*, Comput. Methods Appl. Mech. Engrg., 194 (2005), pp. 1641–1654.
- [9] R. GHANEM AND J. RED-HORSE, *Propagation of uncertainty in complex physical systems using a stochastic finite element approach*, Phys. D, 133 (1999), pp. 137–144.

- [10] R. GHANEM AND P. SPANOS, *Polynomial chaos in stochastic finite elements*, Journal of Applied Mechanics, Transactions of the ASME 57 (1990), pp. 197–202.
- [11] R. GHANEM AND P. D. SPANOS, *Stochastic Finite Elements: A Spectral Approach*, rev. ed., Dover Publications, New York, 2003.
- [12] D. GHOSH AND R. GHANEM, *Stochastic convergence acceleration through basis enrichment of polynomial chaos expansions*, Internat. J. Numer. Methods Engrg., 73 (2008), pp. 162–184.
- [13] E. T. JAYNES, *Information theory and statistical mechanics*, The Physical Review, 106 (4) (1963), pp. 620–630.
- [14] O. M. KNIO AND O. P. L. MAÎTRE, *Uncertainty propagation in cfd using polynomial chaos decomposition*, Fluid Dynam. Res., 38 (2006), pp. 616–640.
- [15] O. LE MAÎTRE AND O. KNIO, *Spectral Methods for Uncertainty Quantification*, Springer, 2010.
- [16] D. LUCOR, C. H. SU, AND G. E. KARNIADAKIS, *Generalized polynomial chaos and random oscillators*, Internat. J. Numer. Methods Engrg., 60 (2004), pp. 571–596.
- [17] Y. M. MARZOUK AND H. N. NAJM, *Dimensionality reduction and polynomial chaos acceleration of bayesian inference in inverse problems*, J. Comput. Phys., 228 (2009), pp. 1862–1902.
- [18] Y. M. MARZOUK, H. N. NAJM, AND L. A. RAHN, *spectral methods for efficient bayesian solution of inverse problems*, J. Comput. Phys., 224 (2007), pp. 560–586.
- [19] H. MATTHIES, *Stochastic finite elements: Computational approaches to stochastic partial differential equations*, Zamm-Zeitschrift für Angewandte Mathematik und Mechanik, 88 (11) (2008), pp. 849–873.
- [20] A. NATAF, *Détermination des distributions de probabilité dont les marges sont données*, Comptes Rendus de l'Académie des Sciences, 225 (1986), pp. 42–43.
- [21] A. NOUY, A. CLEMENT, F. SCHOEFS, AND N. MOES, *An extended stochastic finite element method for solving stochastic partial differential equations on random domains*, Methods Appl. Mech. Engrg., 197 (2008), pp. 4663–4682.
- [22] J. R. RED-HORSE AND A. S. BENJAMIN, *A probabilistic approach to uncertainty quantification with limited information*, Reliability Engrg. System Safety, 85 (2004), pp. 183–190.
- [23] M. ROSENBLATT, *Remarks on a multivariate transformation*, Annals of Mathematical Statistics, 23 (1952), pp. 470–472.
- [24] S. SAKAMOTO AND R. GHANEM, *Polynomial chaos decomposition for the simulation of non-gaussian nonstationary stochastic processes*, J. Engrg. Mechanics, 128 (2002), pp. 190–201.
- [25] S. SARKAR, S. GUPTA, AND I. RYCHLIK, *Wiener chaos expansions for estimating rain-flow fatigue damage in randomly vibrating structures with uncertain parameters*, Probabilistic Engineering Mechanics, 26 (2011), pp. 387–398.
- [26] G. I. SCHUELLER, *On the treatment of uncertainties in structural mechanics and analysis*, Computational and Structures, 85 (2007), pp. 235–243.
- [27] C. SOIZE, *Construction of probability distributions in high dimension using the maximum entropy principle. applications to stochastic processes, random fields and random matrices*, International Journal for Numerical Methods in Engineering, 76(10) (2008), pp. 1583–1611.
- [28] C. SOIZE, *Generalized probabilistic approach of uncertainties in computational dynamics using random matrices and polynomial chaos decompositions*, Internat. J. Numer. Methods Engrg., 81 (2010), pp. 939–970.
- [29] C. SOIZE, *Identification of high-dimension polynomial chaos expansions with random coefficients for non-gaussian tensor-valued random fields using partial and limited experimental data*, Computer Methods in Applied Mechanics and Engineering, 199 (2010), pp. 2150–2164.
- [30] C. SOIZE AND C. DESCIELERS, *Computational aspects for constructing realizations of polynomial chaos in high dimension*, SIAM Journal on Scientific Computing, 32(5) (2010), pp. 2820–2831.
- [31] C. SOIZE AND R. GHANEM, *Physical systems with random uncertainties: Chaos representations with arbitrary probability measure*, SIAM J. Sci. Comput., 26 (2004).
- [32] G. STEFANO, A. NOUY, AND A. CLEMENT, *Identification of random shapes from images through polynomial chaos expansion of random level set functions*, Internat. J. Numer. Methods Engrg., 79 (2009), pp. 127–155.
- [33] N. WIENER, *The homogeneous chaos*, American Journal of Mathematics, 60 (1938), pp. 897–936.
- [34] D. XIU AND G. E. KARNIADAKIS, *The wiener-asky polynomial chaos for stochastic differential equations*, SIAM Journal on Scientific Computing, 24, No. 2 (2002), pp. 619–644.